HIGHLY SCALABLE DATA BALANCED DISTRIBUTED SEARCH
STRUCTURES

By

PADMASHREE KRISHNA

ACKNOWLEDGEMENTS

I am indebted to my Ph.D. committee chairman, Dr. Theodore Johnson, for the

innumerable and invaluable number of hours he has spent in probing me through this

research. There were times when I have come out of his office more confused than

when I went in, but that set me thinking along the right direction. I thank him for

his invaluable suggestions and critical remarks that have led to this research.

I want to thank my other committee members Dr. Earl L. agreeing to be on my

committee and then offering helpful suggestions along the way.

I thank Dr. Randy Chow and Dr. Manjeh Lwidemar my internal committee

member, for steering a lot of interest in my research and for the useful suggestions

that, provided along the way.

Dr. Sterig Reiss offered critical remarks on the research, which led me to re-

think many ideas and things better experiments. Dr. Reiss's more remarks was the

availability of the N-term with large degrees and with a large number of keys.

Lastly, Dr. Richard Newman Wolfe deserves special thanks for introducing me to

my advisor in the first place, providing helpful hints on my research and on writing

the dissertation.

I would like to thank all the CIS staff, in particular Mr. John Bowers, the graduate

secretary, for their help.

# TABLE OF CONTENTS

# LIST OF TABLES

vi

# LIST OF FIGURES

# FULLY SCALABLE FAULT-BALANCED DISTRIBUTED SEARCH STRUCTURES

By

Padmashree Ravindra

May 1995

Present trends in parallel processing and distributed databases necessitate the construction of large volumes of data. Scalable distributed search structures provide the necessary support for such storage. In this research we focus on the performance of two large-scale data-balanced distributed search structures - the dB tree and the dE tree. The dB tree is a distributed B-tree that replicates its interior nodes. The dE tree is a distributed Extendible Hashing that applies lazy updates and uses distributed linear order to represent a distributed index.

The main goal of our research is to develop fully scalable distributed algorithms on large scale data balancing with applications to our distributed search structures. This distributed storage is very safe for data balancing in either the server resource efficiently or with its server resource. The advantage of simulated storage would be bad if there was an application displaced position of the three order with nodes to balance distributed sort large very sequence for lower order to represent a distributed index.

ii

The algorithms for data balancing determine how the best codes are assigned to processors. Here, we develop several algorithms for data balancing, both for the AB case and the AC case. We find that a simple distributed data balancing algorithm works well for the AB case, requiring only a small space and message passing overhead. We compare these algorithms for data balancing in a AC case, and find that the most aggressive of the algorithms tends to keep any data balancing overhead small enough to be worthwhile.

We have also performed some clustering experiments with two different clustering algorithms, to develop some intuition into how clustering affects the overall performance of our system. We find that a simple distributed clustering algorithm works well for the AB case, and that for clustering there are many different subtleties that can introduce overhead in the clustering process.

We have also performed an analytical performance model of the AB case and the AC case using the data from the simulation experiments. We then applied the model to two experimental parameters to obtain an understanding of how the overall data balancing algorithm scales with the number of processors. Our experiments give a value of 40 milliseconds response time with 2 processors, whereas the model predicts 63 milliseconds. From this analytical model, we obtain that a distributed search construct provides a much larger throughput than a centralized index server, at the cost of a modestly increased response time.

# CHAPTER 1
## INTRODUCTION

### 1.1. Objective

The main objective of this research was to develop distributed algorithms and protocols for meta-specific data structures, and implement them to study their feasibility and performance. We approached this problem in two dimensions.

1. Algorithmic: The algorithmic approach attempts to design distributed algorithms for specific sequences and to study their correctness criteria. Distributed data structures are useful to designing general-purpose distributed algorithms. However not all algorithms designed can be implemented efficiently. We are excited to enjoy the implementations of dynamic algorithms on a network of processors without shared memory.

2. Implementation: From the implementation view-point, we were concerned with the efficiency and performance of the algorithms. The implementations of these distributed data structures will help form the application's user, the details of the time where data are stored, the access methods and the synchronization techniques.

For the purposes of this research, we selected the B-tree for its flexibility and its practical use in ordinary large amounts of data.

1

### 1.4.1 Why Distributed Search Structures Were Chosen

Current commercial and scientific database systems deal with vast amounts of data. Due to the volume of data to be handled it is large, it may not be possible to store all the data in one place. Also, often allowing large volumes of data, there is the danger of memory bottlenecks. Therefore, distributed techniques are necessary to create large-scale, efficient, distributed storage [?]. Distributed data structures allow the large amounts of data to be manipulated. The data can be stored by partitioning their among the storage sites of the system, which also allows for parallel access to the data. Distributed data structures are useful for many distributed applications (e.g. in computer information storage and retrieval techniques, global name servers in networks, resource allocation, etc.). Although a considerable amount of research has been done on developing parallel search structures on shared memory multiprocessors, little has been done on the development of search structures for distributed memory systems. One such search structure is the B-tree. The B-tree was selected because of its flexibility and its potential as an indexing large amounts of data.

A distributed system is a collection of autonomous memory-processor nodes in an interconnection network. Distributed systems have several advantages over conventional systems because they enable ease of expansion, provide increased reliability, allow actual geographic distribution, and have a higher potential for fault-tolerance and performance due to the multiplicity of resources. Each processor-memory pair may henceforth be called a site. Sites communicate by message passing. If a balanced tree storage/searching structure partitioned over sites is highly scalable, it is distributed among the sites so that enquiries, insertion and update operations of the global information of the global state of the system. If no such site must have the responsibility of handling accurate and out-of-date information. Distributed algorithms must tolerate this non-homogeneity of

### 1.1.2  The Need for Distributed Data Structures

This data structures used in an algorithm have a considerable effect on the efficiency of the algorithm. Hence for distributed algorithms, there is a need for determining the data structures as follows:

1. The primary reason for distributed data structures is that in a distributed system we wish to share the data between processors on different processors. The various parts of the distributed system share data by communications. Several programming languages only support shared variables (but often for pseudo-parallelism of the processors running on the same processor. Simple shared variables can be implemented by simulating shared physical memory, but this is not sufficient for distributed systems that call for complex data structures. Instead, there are basically three ways of providing the notion of shared data in a distributed system:

   (a) Distributed data structures

   (b) Shared logical variables

   (c) Distributed objects in distributed shared memory.

2. A secondary reason for distributed data structures is the problem of concurrency (or) data structures in one physical location. This not only removes a large amount of contention but also makes the system less fault tolerant. Distributing a data structure over a large number of processors requires partitioning the data structure into parts that can individually managed by a single processor. The parts may be duplicated so they may be replicated to provide faster

faculty and increase availability. But, as explained here is the problem of re-
constraint. Distributing the shared data over the different processors improves
performance since data residing at different sites can be accessed in parallel.

Several programming languages provide the above notion of shared data, where
the user is unaware of the physical distribution of the data.

### 1.1.2 The Principle of Data Distribution

Data organisation must be based on the principle of property defined as:

- all data objects are accessible in all sites;
- in no access the most recent version of the data is provided;
- consistency is maintained on a global basis.

To achieve these criteria, data must be replicated and updates must be ensure.
All these criteria improve the performance and reduce the cost of access by allowing
data depending on the locality of a process. In most data structures, the access index
is predictable, while in others not to.

Data structures are characterised by the operations they support. A distributed
data structure consists of a set of local data structures storing the data at various
sites of the system and a set of protocols for access to the distributed data struc-
tures. These protocols specify the entry and update operations on the data. The
distribution of the data structure is known as the data organisation scheme [?] and
may be based on several criteria as follows:

1. To improve the locality of the process running on the processor

2. To reduce the average complexity of access to remote data

2. To balance the data center processes for efficient usage.

3. To remove the fault tolerance and increase availability.

4. To minimize the delay performance per access.

Resource-constrained on-using processes can elicit the advantage of distributed systems. A good strategy would be to take into account the integration of our communication cost imposed by the underlying machine architecture. Several attempts have been proposed for efficient allocation of data structures [6]. The newer protocols specify the previous operations that are to be performed on the data structures and the mode of access by processes to the data.

## 1.1.4  The Need by Application

Redundancy or replication is an inherent part of the design of distributed data structures. Not only does replication provide fault tolerance in the event of the failure of a processor, but it also enables dynamic data balancing and reduces costs by placing the most often accessed data close to the processor. A process can take advantage of its locality to reduce the cost of communication. Replication also increases the availability of data. A factor that has to be considered is the degree of replication, also known as replication control. In what is called the total structure, all the data are replicated in each processor [48]. This increases the availability and fault-tolerance but places a high demand on memory requirement. A compromise is to set up a balance between memory usage and cost considerations [9]. Replication minimizes the problem of maintaining consistency among the various copies of the data structure.

### 1.1.2  Centralized Data Structure Issues

Interleaving data structures creates new structural parents in a shared memory or a single-processor system. Two basic problems are those created by the concurrency of concurrent data processing operations, and those introduced due to the distribution of the data.

Concurrency issues are resolved by imposing the serializability criteria. Various serializability criteria have been studied extensively in literature in databases [9].

The study of the distribution of data structures and the relationship to the serializing systems of processors may lead to efficient schemes for distributing the data in terms of space, time and average complexity [3]. The complexity of data movements is also an issue for distributed structures.

The much attention releases for the research in the B-tree. We address all the above issues with respect to distributing a B-tree. We have selected the B-tree because of its flexibility and its potential use in defining large amounts of data.

### 1.2  Background

#### 1.2.1  Introduction

In this section, we present a survey on the research done in distributed data structures. Techniques that were proposed to make distributed data structures concurrent are presented. A brief discussion of how distributed data structures is then presented. In the discussion of search structures that the research concentrates on we focus on both tables, dictionaries and containers B-trees. Some background on this technology is also presented. Finally, the concurrent B-tree list algorithms is presented, which forms the basis for the distributed B-tree algorithms.

Distributed memory machines are much more difficult to devise algorithms for than the shared memory machines. This is due in the lack of a single global address space. The programmer is responsible for distributing code and data in addition to that and managing communications between processes. This may reduce programmer productivity, therefore programming languages need to provide facility for developing parallel and distributed programs. In the current conventional programming languages, user process can only access to local address space which results in large data structures that must be partitioned across the processors. Every computation in a computation is usually more expensive than computations; it is essential that track of the computation be done using local data. Several programming languages are being developed to support distributed data structures. Some examples are Linda [1,29], Orca [3, 4] and Kali [32]. Some programming languages provide distributed data structures explicitly when others do so explicitly. Examine the following three.

### 1 Linda

The distributed data structure paradigm was first introduced in the language Linda, implemented in AT&T and Intel's first and developers [1]. The Tuple space concept is used by implementing individual data structures. The tuple space, consisting of tuples that are an an ordered sequence of values. Invoke a global memory shared by all the processes in the system. In namely a out, a "read, modify and write" shares operation is needed. If one processes want to access a tuple, only one of them can read what the other blocks. A distributed array is implemented in a distributed share of memory by each the following systems:

(a) Define the term tuple space is explicitly or

(b) The last processor (or create a tuple at the center of the tuple, or,

(c) A leading function is used to distribute the tuples.

Communication through distributed data-structures is anonymous (as opposed to acceptance communications). Communication primitives such as message passing and remote procedure calls are considered using the tuple space. The processes cannot only through the tuple space. The goal of Linda is to relieve the programmer from the task of parallel programming.

## 2 Orca

The programming language is tightly oriented for developing parallel operations for distributed systems. The data structures are encapsulated into passive objects and can be shared by different processes. The objects are replicated on all processors and are updated by a reliable ordered broadcast primitive [2], [4].

## 3 Kali

A programming environment, Kali is designed to aid in the programming of distributed memory architectures [10]. It allows the programmer to treat the distributed data structures as single objects. A software layer supports a global name space. Algorithms can be specified in a high-level and machine-independent form and the software layer converts it into code of tasks that execute by "message passing". Thus, the programmer is relieved of the task of programming with low level message passing primitives and can concentrate on just algorithm development. The only data type supported is that a distributed array.

### 1.5.2 Distributed Data Structures

Here, we present mechanisms of a ... novel of the basic and ... used style distributed data structures. Scalar variables are usually replicated on each processor.

#### Arrays

Primarily, only distributed arrays are predicted by Kali. Moreover, Kali supports user-defined distributions. Array distributions are specified by a distribution clause [20]. The clause specifies a set of distribution patterns for each dimension of the array. As controls in the dimension indicates on distribution. The number of array dimensions that are distributed cannot exceed the number of processors in the partition. Each processor stores a single copy of each array element.

Another convenient data partitioning scheme has been proposed for distributed arrays [10]. This is a redistribution based approach, whereas the complete analysis each loop nest, based on performance considerations, identifies some constraints on the distribution of data structures. Finally, the compiler tries to combine constraints for each data structure so that the overall execution time of the program is minimized. The data may have to be repartitioned between program segments and between procedure calls. This has been implemented in the Intel iPSC/2 hypercube.

#### Queues

A queue is a First In, First Out (FIFO) structure that has two ends, a front and a rear. A queue can be stored in a distributed system by storing different segments of the queue in different ends, with each queue element being moved to exactly one site.

Level of ... however a distributed fault tolerant distributed queue to provide a high degree of reliability against breakdowns and loss across sites [16]. In this scheme,

... analyses of the open-air stands and each realises a broken near... but necessarily as transformed argument. Each one maximises the local and none of a segment of a region. There are no customary implemented institutions.[20]

### 2.3.4 Stand Structure

Different search structures are needed for establishing lists and abilities as current listed options which have a small primary memory and a larger secondary memory. To access individual address of a file or archive is proposed. The normal operations carried out on an entity are search, insert, and delete. A search table is a data structure in which search are organised in a well-defined manner. Search structures are used for the implementation of dictionaries. An implementation of a search table could be designed using either a tree, an array or a hash table as a sequential structure.

In a traditional design, to access value a long name or complete, usually on the order of the number of elements stored. In sequential-access data structures such as linear sorted arrays and hash tables have need to implement search tables. Of these, the hash table gives the best performance with little space overhead. For address is achieved by applying the account however the sequential nature of the memory renders a bottleneck. Therefore the random access of data structures memory access concurrently is necessary.

Sorted-hash memory data structures have been proposed by Ellis [16], Severance [14], Foley [10], Culbreath et al. [3] and Johnson and Culbreath [30]. Culbreath et al. [3] have proposed a pipelined distributed B-tree, where each level of the tree is constructed in a different processor. The parallelism achieved in limited by the...

length of the $\ell$-tree and the processors are not data balanced. Parallel $\ell$-trees using multi-version memory have been proposed by Wang and Weihl [96]. The algorithm uses a special form of software-controlled cache coherence.

## Hash Tables

Hashing is a well known technique for fast access to records in a large database. One of the main goals of a parallel database management system is to derive several methods of hashing have been proposed which include distributed linear hashing [5], extendible hashing [19], and linear hashing [38] and linear hashing for distributed files [51].

### Distributed Linear Hashing

In linear hashing, the table is gradually expanded by splitting the buckets until the table has doubled its size. Splitting means rehashing of a bucket $h$ and its residents in order to distribute the keys in them between $h$ and one other bucket. Linear hashing requires the use of a series of hashing functions, a new one among every one the table is doubled.

A traditional linear hashing method particularly useful for main memory databases has been discussed by Severance and Fromanek [76]. In linear hashing, the records are distributed into buckets, and each bucket, when its threshold is reached, splits into two buckets. Peak to bucket is hashed, and the data are inserted or retrieved from the appropriate location. For a bucket is located and resident, the record class in the bucket can be placed in any memory module. An index is used to point to the bucket downstream and to control at each processor. Addition computation of the bucket addresses at each processor is small and can in most cases be done in parallel. The index is pointed at the end of data or lowest indexing becomes much address for the hash addressing scheme. When a key and data is entered into the bucket, missing bucket addresses at each processor.

The load copies may be out of data at times causing incorrect hash address

computation. Every hop is used to solve this problem. The paper does of discuss the problem of maintaining local copies of the centralised metadata and recovery mechanism. The design has been implemented on a DBN directly and has some system.

## Extensible Hashing

Extensible hashing combines index search trees (or trie) and hashing. To represent the entry, an extensible trie structure is used instead of a binary tree. The units catch various $2^d$ positions, where $d$ is the number of bits used but currently being used to address the index table. Initially the table structure only one position, which points to a single bucket at one. When the bucket fills, the table is doubled to size, and a new bucket is allocated. When the table is doubled in size, a new bucket is created and the keys are reorganised.

Bits $[d]$ has proposed a distributed extensible hashing technique. As in a sequential system the local structure consists of two parts: the directory components and the buckets. It is the information provided by the directory that allows the buckets to be distributed to different sites in the distributed system and the directories to be replicated among the sites managed and managed by directory managers. The buckets are talked to other sites when they are through a hash link that allows memory stored in one site and managed by directory memory managers. The Extensity manages a essentially a serve example of handling such split requests. The bucket manager is a first rule process that manages a request or replicas of buckets. An operation request on the local table is sent to any directory manager who, in turn, forwards the request to a bucket manager holding a replica of the required bucket. If a bucket manager receives a request to a bucket manager after getting a directory memory managers and managed by directory memory lookup.

The directory manager in this hash is a single member request. A bucket manager to a directory memory request is also their direct to correct the correct. The directory manager has to propagate the update information to all the other

directory manager therefore one problem is that in future affects the entire system. Fault tolerance capabilities are decreased than involve more messages in the systems.

## • Two Phase Hashing

A new hash algorithm for university parallel systems is proposed by Yen and Bastani [I71]. It is partitial systems, claiming gives the best performance, but in massively parallel processors, this leads to a high communication cost. Lowest packing, however, has a low communication cost.

This algorithm called two-phase hashing, combines the clustering and lowest packing strategies. Here, a hash index table is made consisting of node (the hash table keeps those in the table itself instead of being either than clean nodes. If the number of elements hashed at each entry is lowest, then this final factors of each element can be computed. The first phase computation the starting of the distances that an to be hashed at each entry. From this, the first location of each element is computed. The next phase produces the real hashing where the keys are formatted to the appropriate bucket. Than the hash tables there are then formatted in the existing location of this chase. The chain is then searched. A slight increase of few hours packing algorithm leaves as the hypertube hash algorithm is also decreased. In this algorithm, the hash table is mapped directly to the processor space (i.e. the nth entry is mapped to processor 1. Collisions are resolved by rehashing. The difference between the above two algorithms is the method of comparison of the reformed location.

## • Trie Hashing

transmit, uncorrelated subced data structures such as it lines are suitable. However, compressionless DDCS provide for dynamic flow we indeed explore. The fundamental algorithms builds the file such the more top space as a fl you but without its nodes by using col.band. Two other algorithms we have throughout of the interest by adding the orders so other the circuits, as we line chains and the errors simultaneously reducing the truck test.

Distributed file compression for such nuodes files has been discussed by Vangourk et al. [20]. The focus at their work has been to achieve scalability (in terms of the number of servers) of the throughput and the file rate while dynamically distributing data. These results indicate that scalability is achieved in a controlled such performance.

## Reference

A dictionary is a dynamic data structure that supports the following operations: insert, delete and search (lookup). A search at the most fundamental data structure and is useful in many applications, such as natural language systems and database systems, in maintaining symbol tables and picture matching systems. In a text-based system, operations on dictionary is a function of the number of elements. Exploiting the operation will give more production, but this task is to be limited to the more frequently accessed cross. The bottleneck between errors as the number of processor increase is a simple alternative to worth supplemented in two trivariants such as the AES case and the BT case. The emphasis of the sequential dictionary structure is a logarithmic function of the number of elements.

A sequential dictionary structure has been proposed three simultaneous and rebalance structures [33]. The objective of this system is to remove the sequential

concentration. The risks graph is similar to the hobs such as a 5-min DP and allows fast sequential access. The index or the abbey graph are used to traverse the entire structure, hence the duration of the graph must be fast and duration of the graph points on angular graph.

Poly [38] has presented a detailed example of a distributed data structure. A compact dictionary structure called BPL is described. The BPL is based on a "flat" tree consisting of an index - a sorted vector serving as the directory and a collection of *leaf* (each maintained at some index) that make its data as an orderly fashion. The paper also discusses the distribution of the central server and replication issues. Complexity issues and memory behaviour are also addressed.

Search structures based on Linear Ordering Linear Structures (LOLS) finally [such as $B^+$ trees, $R/K/S$ trees, etc.] have been proposed [43]. The paper addresses the problem of designing smart structures to fit shared memory multiprocessor and serial search structures. The locality of the structure is partitioned into a number of identical sub-vectors (the sub-vectors have the same structure and randomly) which are stored in the shared memory such that the access they structure over the keyboard over the processors. The design goal is the feature that client data memory independent white being the index of its client and rely throughout the data from being in single same upon servers and being able to spread index function.

## 1.3 Contribution of this Work

This dissertation addresses several issues such as fully distributing a H-tree, its various address and ensuring of a scale data balancing and replication, among others

Data partitioning also raises new issues such as allocating storage for the data, efficiency of access and balancing data among processors. The balance of concern for distributed storage are throughput, scalability and reliability. Most of these issues of interest are available in the current literature; but not in a consistent with each other. Our work addresses all these issues.

We have developed a theoretical framework for explaining the present notion of Live B-tree. Based on that, we have implemented two strategies of replication, namely full replication and path replication. The performance of these algorithms show that the path replication is better and is more scalable. We have developed several algorithms for data, balancing a distributed replicated B-tree. We present the performance of our algorithms. An application of the work to distributed cursor (see, [28 new]). We developed several data balancing algorithms for the distributed cursor too.

1.3.1 Structure of the Dissertation

We have organized this dissertation into four broad categories: theory and practice. In Chapter 2, we provide background on structures B-trees, the distributed B-tree and the distributed cursor too.

Chapter 3 provides the theoretical framework for explaining the notion of a B-tree. In Chapter 4 we present the implementation design details. We present the underlying architecture and message passing mechanism for our implementation. We also present some prototypical protocols that are common in all our data balancing algorithms. Finally, we discuss the portability of our implementation from the SUN to the KSR-1 shared memory parallel machine.

The performance of our replication and data balancing algorithms are presented in Chapter 5. Here, we discuss the replication strategies and discuss the results on their performance. We next discuss the various data balancing algorithms on the

dilctate and compare their performance. We also present the performance of the GliderT.

We conclude the dissertation by summarizing the contributions of our work and providing some ideas about the direction for future research.

# CHAPTER 3
## SURVEY OF RELEVANT WORK

## 3.1 Introduction

In this chapter, we present some background on concurrent B trees, concurrent Δ-link algorithms, the distributed B tree, and then following the distributed B tree. We also provide a discussion of the paper by Johnson and Colbrook [JC]. They introduce a new balanced search tree algorithm for distributed memory systems. They can be. It led me to a basis for the distributed B tree, the dB tree. To reduce the cost of maintenance of the distributed B tree, a path replaced a strategy is used, whereby if a processor owns a leaf node it also owns all the nodes from the root to that leaf. The replication of the root is to every processor renders operations to be actuated at any processor. The leaf level nodes are not replicated. The concept of data balancing has also been introduced to balance the load at all processors. They prevent some sites on base links following can be explained using distributed B tree tree algorithms. Finally, they also show how the dB tree algorithm can be used to build a data-balanced distributed dictionary: the dB tree.

## 3.2 Concurrent B trees

Tree structures (in particular B trees) are suitable for creating indexes. B trees of high index are meanwhile once they result in a reduction of the number of disk accesses needed to search an order. If the index has P entries, then an order to α = k + 1 would have only one level. An structure which causes a node to become too full adds the cost of restructuring of the tree to the overhead of the

Current database designs materialize the construction of databases which allow for consistency of stored processes. The original B-tree algorithms were designed for sequential applications, where only one process accessed and manipulated the B-tree. The main concern of these algorithms was minimizing access latency. However, with the growth of processing power and the need for parallel computing, concurrency throughput has become important. The B-tree is suitable for concurrent operations by allowing multiple processes to perform subsequent operations.

Several approaches to concurrent access of the B-tree have been proposed [5], [18], [34], [22]. All the algorithms share the problem of contention, which can be categorized into two types: data contention and recovery contention. Both lead to performance degradation.

- *Data contention.* All concurrent search tree algorithms require a concurrency control technique to keep two or more processes which access the B-tree from interfering with one another. This contention is most pronounced at the higher levels of the tree. All algorithms proposed use some form of locking technique to ensure exclusive access to a node.

- *Recovery contention.* Performance degradation or unreliable when several processes access a single resource in the system. In shared-memory, concurrency and recovery techniques must ensure correct data when one process wants to access some data and other process wants to modify the same data. A look-ahead mechanism has been proposed to minimize contention where one process accesses data and many other processes want to modify data simultaneously. The support and recovery techniques must ensure correct data after restart, and the data in super processors [26], and Lehman and Yao [18] are a lock technique to reduce contention.

Parallel B-tree using multi-cursor operators have been proposed by Wang and Weihl [40]. The algorithm is designed for software cache management and a suitable

```
/*** B-TREE ***/ node
int64_t java::node::TBcChild(node){
    struct spec node TBcChild(node){
        Bnode   Bnd
        extern     /* nhandle ja a (nulate */
        return     /* search a q nmdlete */

   /* a nilxtecuninvl is   */a

   static protected *a

        struct java::node *node;
        node = (node)
        Bkeyu q hed(a)
        while (node !null)
        {
            child = k.BnodeectruetnsetBe_a)
            if(a < child) child
                node = node[child]
            else if(a == child)
                return node
            Bkeyu q hed
        }
    /* internal a right, a left reason pair p,a , a a a */
        child = BnodeectnsetBe_a)
        int++ Bchild key (a, a++a)
        if(a = BnodeectnsetBe)
            /* a if balia */
            *a
            node = Bchild(p) a
        }
        %p ;a (a
        return (node)
    }
    %a ;a (a
    return (null)
}
```

Figure 3.1: Search Algorithm for a B-tree Tree.

Figure 3.1: Roll-split Operation

hold at most unit can hold at a time, restructuring must be separated into disjoint operations. The first phase is to perform a *roll-split* operation (Figure 3.2). During this phase, a new node (the sibling is created and half the keys from the merged node are transferred into it. The sibling is put into the leaf list and the sibling pointers are adjusted appropriately. The next phase is to inform the parent of its split. Since the lock on the leaf node is released, the parent node is locked, and a pointer to the sibling is inserted into its parent. During this time that the split occurs and the pointer is inserted into the parent, operations navigate to the sibling via the *link* field, not the *highest* field or the node.

On the fly node deletion is not supported in these current configurations. Several dimensions as an site-by-site no exact, including never deleting nodes, performing garbage collection or having the deleted nodes as sinks without disallowing them physically.

## 4.3 The Hitter

Johnson and Collusch [9] present a data-driven hitter-model for average pass-ing architecture. The soccer striker are subjected to in-game parallelism and ulti-mate the landmark. The enormous state sense is Led runs off the strides in the path from the cost in the leaf. Instructing timeshots are made locally, thereby reducing the communication overhead and increasing parallelism. The paper also deals with the time balancing among processors.

The 2K tree is built upon the innermost B-tree algorithms. In the offense, the leaves are distributed among processors. The interior nodes are replicated among the processors. Every processor at a level has links to both its neighbors. Also, each node stores the distance from its leaves. Nodes of the different are given unique keys. A processor maintains a radio contact on the revisions of a node. The tag is a concatenation of a node number at a processor and the processor number. A translation table is used to access a node.

The uppermost most-stable and used are defined on the 2K tree. Concerned rep in such operation, unions are performed on the nodes of this tree. A processor accepts messages from other processors for performing the operations. Messagepool messages are routed to the correct processor. When a node becomes full, a "halfsplit". The inside links of a node help in performing the half-split. Initially when a node merges into another node (becomes empty), a merge is said to happen. The left link stays the same. A full range passenger is sent. All links in a merged node must be changed before a merged node can actually be deleted into the tree.

## 4.4 Architecture

The multi-version memory algorithm presented by Wang and Weihl [10] utilize the amount of non-locations and communications implicit in enactions replicated

so that which must be achieved with all other actions, as their separate communications with other needs.

Johnson and Ayodeya [30] present a framework for creating and analyzing hot market algorithms. The framework is used to develop algorithms that can manage a riff over node. The algorithm uses key smart actions and state synchronized half node actions. In addition, the algorithm framework accounts for cultural actions to ensure that classes of actions are performed as a node as the actions to which they are processed (i.e. the hub charge actions are achieved).

## VI. Data Repository

To avoid substantial storage space utilization at processors, it is necessary to perform data balancing on the processors. The balancing also spreads the caches to the data structure evenly among processors. It also provides rapid recovery and cause utilization at each processor.

Data balancing among processors has been studied by Johnson and Culianush [30]. They suggest a way of achieving communication rate for data balancing by moving neighboring leaves in the same processor. When a processor decides that it has too many leaves or data, it processes it at a processor leaves. If that processor accepts the leaves are considered, the processor will be neighboring processor is highest, thus it balances the data by taking balance the loosely loaded processor leaves for a lightly loaded processor and transfers the lower.

In the context of code modeling, object modeling has been proposed in Kumahili [20]. Algorithm keep balancing subsystems even after they have moved so another node call use a homefrost partition if its forwarding submanian is available.

Lee et al. [20] have discussed a fault tolerant scheme for distributing system. The scheme described by them provides dynamic fault tolerance, high availability and uniform load balancing with small and error space requirements and low communication. High availability is achieved by replication of the queue and task queue replica may be distributed over several sites. Consistency is maintained by two-phase locking. Small storage space is needed in each processor, since only segments of the queue may be kept in a processor. Since global broadcasting is not used, the communication overhead is low. However, since queue access requires communications to access global consistency. When a processor locates a queue repository, it sends a request to the processor containing the head or tail of the queue: the receiving the request, the current head or tail processor will lock up all other head or tail queue replicas thereby ensuring consistency. If the processor which receives a request does not hold any head or tail replica for the segment, it looks for the nearest neighbor holding the head or tail is found.

Effici algorithm [24] performs data balancing whenever a processor runs out of storage. Telag [44] has studied the issue of data balancing in distributed dynamically from a completely new point of view, stating the processor state transition (sleep, busy, blocked) of the number of keys and 'P' as the number of processors. In general, this algorithm mathematically for both sleep and in weak forward algorithm converges and exits repeated.

In the off-line: this data balancing is performed by distributing the keys among the processors. This requires communications among the processors only once a key moves to restore ability and queue holes. Also, the number of internal nodes each merid is high: An alternative to the off-line is an off-line.

## 2.4  AX-tree

To reduce the communication cost, Julienne and Dubhashi suggest the AX-tree, also known as the distributed crystal tree, where neighbouring leaves are stored in the same processor. They define an order to be a maximal length sequence of neighbouring leaves that are stored by the same processor. What a processor decides that it sees has many more: it first looks at the processor who own neighbouring entries. If the neighbour will retain the leaves, the processor transfers some of its leaves to the neighbour. If no neighbouring processor is lightly loaded, the heavily loaded processor searches for a lightly loaded processor and creates a new route.

Figure 2.4 shows a four processor AX-tree that is then balanced using the new tree. The routes have the characteristics of a leaf in the AX-tree: they have no upper and lower range, are doubly linked, access the intervening operations, and are numerically safe or ranged. The initial balanced AX-tree can be treated as a AX-tree. Each processor stores a number of routes. The least stored in the route are kept in a sorted linked list structure. Each route is linked with its neighbouring route.

When a processor receives a search request on a AX-tree. When a processor knows that it is on locally loaded, it looks at the neighbouring entries to take more of its own. When there is the AX-tree. The route corresponds to splitting it by using more operations, or a new AX-tree. This helps reduce by using all neighbouring routes as loosely loaded, or when both, or the same number of routes as to move more to AX-trees of splitting operations.

Each operation can make some keys, the route can move a processed so the number of processors. Also, each restructuring is greatly reduced as it takes place only after a large number of keys have been inserted or deleted.

The AX-tree can be used to address shared file systems (DFS).

Figure 2.2: The 42-Arm

## 2.9.1 Beyond The Systems

Parallel file systems have been proposed to better assist IO throughput in processing power. A parallel filesystem in a file system in which the files are stored on multiple disks and the disk drives are located on different processors. A common method for implementing a parallel filesystem is to use disk striping [51], in which successive blocks or a file are stored on different disk drives, such that two or more successive A parallel striped file system, Design has been implemented in the MPP Intel/ely [43]. A striped file can be expanded (or prepended) to and accessed its structure. However, a block can't be accessed once at defined into the middle of the file; the data stored in routing the regular engine structure of the file, because of this mention of the nature (or regular engine) is the 42. A manipulation of the file is a requested design. however, does not support three operations. In many applications, the most common operations on the IO are "read" and "append", as deserve volume-latency. Certain other applications use "insert" and "delete" from the middle of the file.

In the universal crystal the proposed fee indicates (Fig. 22), the the contains of a single critical axially. An insert to deform calls for decisions to be made about an organising structure. The particles to attend the fact must a split correspond to the ordering of a unit in size. It was used, passing of time extends corresponds to a merging of the strides. Thus, the 13-own dependent size pursue an better structure which allows one to amove axis or slides from a to-ged flu, and further, so the oblegid current are faxed together. the flu can so-sequentially used as a higher-parallel current to provide last renden motor. Insert more to the ten.

The conception is what the flu is composed of crystals, each of which can be atterated by a flag which in turn can be ordered. The assumption to crystallite because the crossing of "insert flu thus slide, on the fau" and the "fau" the proveny means date break are overlaid, and on a higher-parallel current to provide last renden motor. Insert more is the ten.

The conception is what the flu is composed of crystals, each of which can be atterated by a flag which in turn can be ordered. The assumption to crystallite because the crossing of "insert flu thus slide, on the fau" and the "fau" the proveny means date break are overlaid, and on a higher-parallel current as a given value of a possible attered the of the connection of the a lower structure date because the motor is not ordered. The assumption is crystallite become the crossing of "insert flu thus slide, on the fau" and the "fau" the crossing into the ordered at a new motor control. The the also because it has no to tetrome, within the rotor the its integround or a new motor control. The the also acts helps to emerge the doped return.

An example of an ordered output flu is shown in figure 2.4. The flu is broken into number of strides, each of which is comparatively a used across at basis (i.e. a single axial). The onion are used according to the adopting the structure, which will be to-day in the resulting to order the renden motor.

Figure 3.1 An Isolated Helped Film

## 3.2. Conclusion

In this chapter, we have presented a background on concurrent B-tree and the *standalone* B-tree. We have discussed the work done by Johnson and Cotilhard [20]. They pointed some ideas on the implementation of a *standalone* B-tree and also present some techniques to avoid bit-rest bad caused by replication of the advisor index. Further, they discuss some ideas on balancing the persistent single field the distributed B-tree. Concurrency control and replica coherency are also addressed. To extent the cost of centralization for data balancing, they support the distributed action tree. The diff-tree manages certain extend of contextual keys. Johnson [27] provides a discussion of how the diff-tree can be used for a practical application of staged the options. In the next chapter, we provide a theoretical framework of this algorithms for replication of the *standalone* B-tree.

# CHAPTER 2

# REPLICATION ALGORITHMS

## 2.1 Introduction

When addressing large volumes of data, there is a danger of memory bottlenecks, where all processors access the same data item stored at one processor. For example, one of the problems with a distributed search situation is that since all accesses to the data have to pass via the root node, the root node becomes a bottleneck and correspondingly node which stands at the top of the hierarchy. It also causes excessive message traffic as the network transports the processes which holds the root node of the search structure. This is known as a memory contention and can be solved by replication. Allowing multiple copies of often accessed same databases can be made least among the components of the system. Replication while providing redundancy, reliability and improving contention, however, introduces consistency problems previously not present. A method of achieving consistency is to guarantee that all operations take place in the same order as did the operations of the distributed system.

Several algorithms have been proposed for maintaining a node [?]. These, however, do so at the cost of concurrency since they require synchronization and thus entail significant communication overhead. Lazy replication has been proposed by Ladin and Liskov for enforcing serial ordering. The issue seems to be loosely completely serialized, in spite of the ... causal relation. Explicit ... based on ... through ... The cost seems ... by introducing explicit dependency among messages. This, however, comes at the cost ... they are causally related but enforcing in arbitrary replicas. A set B of ... it to imposition on the previous ...

ual, a, the region which receives δ, or, δ, does not have enough admittance about it to proceed. The replica, δ, has to δ(for) the insertion of δ and δ λ receives all the updates δ λ(patch) as...

Performance cost to reduce the rate of maintaining replicas/data and for accuracy and consistency. Later, Ladin, and Shira propose lazy replication by maintaining replicated services [SK]. Lazy replications are the dependencies that exist in the operation issues to determine if a service's data is sufficiently up to date to maintain a new request. Several authors have explored the maintenance of new checking and note free context and data structures in a shared memory environment [11]. These algorithms estimate consistency because a slow operation event blocks a fast operation.

In this chapter, we present an approach to maintaining distributed data attributes whose attributes may be updated, which take advantage of the semantics of the read-direction operations to allow for variably-sized but localized replicas. Lazy updates can be used to change database event dependencies that support the level of consistency. The alternative to lazy update algorithms (represent updates) use synchronization to ensure consistency.

Lazy update algorithms are similar to lazy replication algorithms because both are the mechanism of an operation to reduce the cost of maintaining replicated copies. The effect of an operation can be lazy sent to the other servers, perhaps as part of each subsequent messages. The lazy replication algorithm blocks an operation until the local data is sufficiently up-to-date. In contrast, a non-blocking read first processes data structures never blocks its operation. The lazy update algorithm we propose so that the execution of a recent operation never blocks a local operation; hence they are a directional technique of non-blocking operations.

Lazy updates have a number of properties advantages over more expensive replicated operations. They significantly reduce synchronization overhead. They are highly...

Figure 3.1: A file tree

### 3.3. Replication

All questions start by assuming the root of the search structure. If there exist any one copy of the root, then access to the index is contained. Therefore, we want to replicate the root widely in order to improve parallelism. As to increase the degree of replication however, the cost of maintaining coherent copies of a node increases. Since the root is vastly updated, a significant portion of the accesses are updates. As a result, wide replication of leaf nodes is productively expensive.

In the dB tree the leaf nodes are stored on a single processor. We apply the rule that if a processor stores a leaf node, it stores every node on the path from the root to the leaf. An example of a dB tree which uses this replication policy is shown in Figure 3.1. The different replication policy causes the cost tradeoffs, the lower as a single processor, cost be to processed to accommodate. The root is replicated as a maximum to provide maximum parallelism. As a result, updates to the root are infrequent relative to the mounting traffic, reducing the number of messages passed.

The replication strategy for a dB-tree helps to reduce the cost of maintaining a distributed search structure, but the replication strategy alone is not enough. To very work within imposed the accesses of to available copies algorithm [3], are involved

Figure 3.2: Lazy scans

of maintaining replicated copies would be prohibitive. Instead, we take advantage of the semantics of the extent on the search structure nodes and use lazy updates in situations like replicated copies propagated only

We note that many of the schemes on a different node constraint. For example, consider the example of schemes which occurs in Figure 3.3. Suppose that node A and B split at "when the scan time." Features within one sibling must be counted into the scan. In this case, the route time needs a square. A pointer to A' is inserted into the search structure, and the first copy of the parent with a pointer to B' is carried into the search structure. Now we copy the parent and continue to concurrently scan not only does the parent not contain a pointer to one of its children, but the first copy of the parent don't contain the same value.

The scan in Figure 3.2 is still works, since we only has been made unreliable. Further, this report of the priority will eventually converge to the same value. Since both lazy, since as soon as the scan will eventually emerge to a node the other node. Given that, as long as the priority will eventually converge to a node in the synchronism with respect to updates in the scans, it doesn't need a fix. This is important because we maintains an inversion with these later synchronization requirements key updates. We still note concurrent with root there synchronization requirements key updates.

## 2.2 Overview of Distributed Search Structures

Shasha and Goodman [9] provide a framework for proving the correctness of two replicated concurrent data structures. We make extensive use of their framework in order to discuss operation correctness. We detail most details here in two spaces, but we note that if the distributed analogue of a "scalable-replica" algorithm [25] illustrates we may not be less of the distributed concurrency. In this section, we discuss a model of distributed search structure computation and establish correctness criteria for lazy updates.

A node of the logical search structure might be stored at several different processors. We say that the physically stored replicas of the logical node are copies of the logical node. We denote by copies $(v)$ the set of copies that correspond to each a of (global situation) node $v$.

An operation is performed by traversing a sequence of actions on the copies of the nodes of the search structure. Thus, the specification of an action on a copy of the node involves (perhaps via an "intermediate action" algorithm [25]) the ...

A first value $f$ and a subsequent action set $S_f$. As soon as the specification enables a node of $m$ copies stored at a processor, it defines the actions that can be performed at a value before the operation that is performed at one of the copies that ...

adapted to the remaining inputs. We distinguish between the initial inputs and the original inputs. Thus the specification of an action is

$$\alpha(p, c) = (c', S)$$

When action $\alpha$ with parameters $x$ is performed on copy $c$, copy $c$ is replaced by $c'$ and the subsequent actions on $A0$ are scheduled for execution. Each subsequent action is of the form $(\alpha, p, c)$, reflecting that action $\alpha$, with parameters $p$, should be performed on copy $c$. If copy $c$ is scheduled, the parameter $p$ for action $\alpha$ may be performed on the copy $c$, if it is an actual copy. Otherwise, copy $c$ reflects the copy that could not be ready because it is still waiting for an actual copy $c$, and $c$ can be an actual action by writing from is internal ...

In order to discuss the consequences of writing, we will need to specify whether the order of two actions can be exchanged. We say that two actions $\alpha$ and $\beta$ commute if performing them in either order results in the same result, and their effects are independent. Two actions that do not commute are said to conflict. If the two actions $\alpha$ and $\beta$ both write to ...

In order to discuss the consequences of writing, we will consider ...

An algorithm might require that some actions must be performed on all copies of a unit, or some other kind of "transformation." These serve some other effect, which is essential to the extension of $AA0$. The outcome of an $AA0$ operation is an $A00$ ... A copy may use or modify $AA0$ recursively ... and conflict with others. We assume that ...

the node manages at each processor is aware of the AAA strain conflict relationships, and will block strains that conflict with currently executing AAA. The AAA on the downstream processor manages in the following way.

### 5.1.1 Examples

In order to capture the consistence nodes which actions on a copy constraint, we model the value of a copy by its factory in $n$ [10]. Formally, the total factory of copy $v$ is represented amount of the pair $(I_v, R_v)$, where $I_v$ is the initial value of a copy and $R_v$ initial value of an action $n$. We define successions in terms of the $R_v$ so factors. The initial value of the update action $u$ ... We denote the succession between amount and new object actions, we define the succession between $X[u]$ to be the opinion factory. If each such action $u$ realized $u$ with each amount of realized in $u$. Finally, we write over the initial factory of the latter $i$ so $u$ with amount $n$.

Because that $R_v = \Delta R_{v+1}$ and this $\Delta$ the first value of $R^* = I^*_{v+1}$. Then $R^*_v = (I^*_v, r^*_v)$. $(R^*)$ so with the factoristic situation of $X_v$ by $R^*$ in a way to each $w \in R_v$ over $w$ the but thus the factory also. Finally, a also is consider. For a natural value, $I_v$. When a copy of a node at created it is given an initial value, which will fact amount over $R_v$. The initial value of ... we call the composed value of the copy. That initial value consists the action summarized over the factory of copy $v$ by an natural value such of the copy and will typically be equivalent to the factory of the executing node $n$ at a the situation of the factorists of the executing node to the factory of the initial of update node at a the matter of ... over the factory of update node matters of ...

performed in the copy. If a copy of a node is deleted, then we no longer need to worry about the node contents. We denote a set of all initial updates not yet performed on node $n$ by $m_n$.

We recall that an action on a copy or value if the action on the primary value of the copy has an associated subsequent actions. A history is valid if a node is modified as a copy of the primary value of the node, along with all initial updates not yet performed on the node, is applied. We denote by $E_i$ the set of values not yet performed on copy $c_i$. Let $R_i$ and $W_i$ be the set of read and write actions on the copy $c_i$. Let $A_i$ and $N_i$ are compatible, then we write $A_i A N_i$.

Our correctness criteria for the replica maintenance algorithms are the following:

**Complete History Requirement.** A node is with initial value $k$, and update actions on $M_n$ has compatible histories if, at the end of the computation of $C$:

1. every copy $c$ [unreadable] with history $H_c$ has a landmark metaevent $d_c$ such that the action actions on $M_c = E_c D_c$ contain exactly the actions in $M_n$,

2. every landmark metaevent $D_c$ not be removed if it is from $M_n$ and that $D_c[H_c] = N_c[H_c]$ for every $c_i$ and every $N_i \in M_n$.

If an algorithm guarantees that every node has a compatible history, then it meets the compatible history requirement.

**Complete History Requirement:** If every subsequent action must appear in some node's update action set, then the requirement meets the complete history requirement. Conversely, assuming completeness the algorithm maintains set for the complete history requirement.

**Ordered History Requirement:** We define an ordered action as one that belongs to class $c$ such that all actions of class $c$ are compatible with each class (we assume a total order exists). A history $H$ is an ordered history if for any ordered

nclme $k_n, k_t \in R$ of class t, $\partial k_t < \partial_t$ then $k_t < k_0$ in $R$. An algorithm enters the *colored history* represented if every node has a compatible history that is in *colored history*.

The reversible history requirement guarantees that every node a single-step equivalent when the constructive recequences. We note that our condition for every property uniform history as a condition of the subsequent action any rather than a condition of the intermediate values of its state. The reason until only to have the same value at the end of the transposition, but the subsequent action can't be problematically stored as subfactors without a special protocol.

The complete history requirement tells us that we must touch every record visited in a copy. A default node is conceptually situated in the search situation to satisfy the complete history requirement. The colored history requirement lets us recover *explicit synchronization requirement* on the equivalent, parallel algorithm by evoking the conditions in the copy inference algorithm.

### 3.4.2 Lazy Updates

An update action must be performed on all copies of a node. With an further information about the action, it must be performed on a full in reason that. So invalidating actions are ordered in the same way with all copies. However, some actions associate with class a *lazy update* attribute, when acquiring the need for an AM. In figure 3.1 the final value of the node $a$ is the same at either copy, and the evaluation is always at a good node. Therefore, there is no need to agree on the order of execution. We provide a simple description of the *degree of synchronization* different *update* acquired.

**Lazy Update.** We say that a *much* adoption update or a *lazy update* if it associates with all other *lazy update*, as synchronization is not required.

Down-regulatment option. Other markets are almost key options, but they conflict with some other actions. For example the actions may belong to a class of ordered actions. We call these meta-synchronous options. A meta-synchronous action requires special treatment, but these are outside the estimation of an AMI.

Synchronous Update. A synchronous update requires an AMI for correct execution. We note that the AMI might block only a subclass of other actions, or might extend to the region of two (or) different nodes.

## 3.3. Algorithms

In this section, we describe algorithms for the key measurement of several different AMIs (or) algorithms. We track from a simple fixed-expert distributed theory to a more complex credible expert X tree and develop the algorithms we need to deal with along this way for all of the algorithms we develop, we assume that only search and action operations are performed on the AMI tree. In addition, we assume this network is reliable, delivering every message exactly once in order.

### 3.3.1. Fixed-Partition Layout

For this algorithm, we assume every node has a fixed set of experts. This meaning this lets us concentrate on specifying key options. Every node contains pointers to its children, its parent, and its siblings. When a node is created, an ordered set of experts is also made present. The node-id pointer, the root expert entry for its parent (and) its sibling nodes, this pointers if this tree and becomes too full, the operation

corresponds the fill text by arrows half-right and insert actions. The insert action adds a new key to the known and adds a pointer to a child as the new leaf nodes. The half-split action creates a new sibling (and the sibling's copied), transfers keys from the full node to the sibling, corresponding nodes to point to the sibling, and adds to an insert action to the parent, so its parent in the internal node.

The first step in designing a data-driven algorithm is to specify the canonical storing relationship between actions.

1. Any two insert actions on a copy commute. As in Inger's algorithm [69], we need to take care to perform out-of-order inserts properly.

2. Half-split operations do not commute. Insert a half-split action modifies the right-sibling pointer, the final value of a copy depends on the order at which the half-split are processed

3. Delayed half-split actions commute with adjacent inserts, but not with undelayed insert actions. Because this is Inger's $P_2$, range Inger's. Thus, if the order of $i$ and $x$ are switched, $x$ becomes an invalid action. A delayed insert action has no subsequent inserts, and the final value of the node is the same as the empty or delete ordering. Therefore, delayed half-split and adjacent inserts commute

4. Initial half-split actions don't commute with adjacent insert actions. One of the subsequent actions of an initial half-split action is to create the new sibling as ordinary. The key is to send a delete action to commute as possibly half-split node to which the subsequent action an inline insert at point the half node. The subsequent in either a insert or where the sibling depending on whether a insert before or after the half-split, depending on whether a insert before or after the half-split.

By our classification methods, an octet is a *lazy* update and a *half* split is a *semisynchronous* update. If the ordering between *half* splits and octets is/are exact (recent the result is *best* update (see Figure 3.3). We next present two algorithms to manage fixed copy octet. In order to list split fields, a computation use a process copy [FS], which executes all set to *half* copy actions: (one *FC* copies never mount recent *half* split octets, only relayed *half* split). The algorithms differ in how the recent and *half* split actions are relayed.

The synchronous algorithm uses list the octet and *half* split as a recent action. The synchronous algorithm uses the set to the square root (slower). The non-synchronous algorithm requires that the ordering of the primary copy be consistent with the ordering in all other nodes (see Figure 3.4).

We do not require that all fixed octet actions are performed at the FC, or require might but that they exceed their recurrent capacity. However, once each copy is executed actually, it is a simple matter to add another blocks.



Figure 3.3  An example of the hot-octet problem.

## Synchronous Rules

**Algorithm.** An operation is executed by subscribing to action, and each action generates subsequent actions until the operation is completed. An operation is executed by executing its B-list (set action) as described previously. Thus, all we need to do is specify the execution of an action at a copy. The synchronous split algorithm

uses an AAE to ensure that rights and secrets are entered the same way at the PC and at the user PC input (see Figure 5.4).

*Half-split* Only the PC executes initial half-split actions. The PC copies secrets relayed half-split actions. When the PC detects four or more half-split the make, it does the following:

1. Performs a *calculated* AAE locally. This AAE blocks all stored *secret* to make, but not relayed based on stored *secrets*.

2. The PC sends a split *start* AAE to all of the other *copies*.

3. The PC waits for acknowledgement from all of the copies of the AAE.

4. When the PC receives all of the acknowledgements, it performs the *half-split*, creating all copies of the new *clicking* and *sending* them (for relay's replied value.

5. The PC sends a *split* AAE to all copies, and performs a split *end* AAE on *itself*.

When a user PC copy *receives* split *start* AAE, it blocks the execution of *initial* events and *sends* an acknowledgement to the PC. The execution of broke *initial* events but *waits* on acknowledgement to the PC. The execution of *initial* events *resume* when the user PC copy *receives* the *split* AAE, it *unblocks*, performs the *copies* and AAE *on itself*.

*Insert* When a copy *receives* an *initial insert* action it does the following:

1. Checks to see if the *insert* is at the copy's *range*. If not, the *insert* action is *sent* to the *right* *sibling*.

2. If the merit is in range, and the step is performing a split AAB the merit is blocked, otherwise

3. The merit is performed and infeasil merit actions are sent to all of the other cases.

When a copy receives a relayed merit action, it checks to see if the merit is in the copy's range. If so, the copy performs the merit. Otherwise, the action is discarded.

Result When a copy receives a merit action, it consumes the unit's current state and starts the appropriate subsequent action.

We saw that once non-FC copies can't achieve a half split action, they may be required to perform an assist or a can-half node. Assists are a way performed on a single process, or it is sent a condition in which a temporary overflow bucket. The FC will soon detect the overflow condition and raise a half split, initiating the problem.

**Theorem 1** The asynchronous split algorithm satisfies the complete, nonsplit and ordered history requirements.

**Proof** We observe that the fourth FIA algorithm gradient is satisfied, as that when ever an action arrives at a copy, it's guarantee to satisfy the copy's nonsplit. Therefore, the synchronous split algorithm satisfies the complete history requirement.

Since there are no ordered actions, the synchronous split algorithm necessarily satisfies the ordered history requirement.

We show that the synchronous algorithm produces compatible histories by showing that the histories at each node are compatible with the actions history in the

partway cage. First, consider the ordering of the half-side actions in half-side is performed as a node when the quiz and AAB is rewritten. All initial half-side actions are maintained in the PC, then are relayed to the other output cages. Since we assume that messages are received in the same start, all half-side are presented in the same order at all nodes.

Consider an ordered event $I$ and a coupled half-side $y$ rewritten at one PC type $c$. If $I < c$ at $D$, then $I$ must have been performed at $c$ before the AAB start for a coupled at $c$ because the AAB start must. Therefore, $D$'s relayed must $c$ must have been sent in the PC before the acknowledgment of $y$ was sent. By message ordering, $c$ is ensured at the PC before $I$ is performed at the PC, so $c < D$ at $P$. Thus, if $c < D$ at $D$, then $c$ is [also] at $P$; if $I < c$ at $D$, then $I$ is [also] at $P$; else $I = c$ (due to message passing causality). $\square$

We note that this algorithm makes good use of type updates. For example, only the PC sends an acknowledgment of the split AAB. If every channel of concurrent tree between cages had to be locked, a split action would require $O(\mathrm{connect}\cdot\mathrm{C})$ messages instead of the $O(\mathrm{cages}\cdot\mathrm{C})$ required for this algorithm. Furthermore, much actions are never blocked.

### Same synchronous hiatus

We can greatly improve on the synchronous-safe algorithm. For example, the synchronous-safe algorithm blocks until events when a split is being performed. Furthermore, it is [requisitely] messages are required to inform the exit. By applying the "real" of ensuring access ordering, we can obtain a simpler algorithm which gives the benefits of [sequence-safe] without its [drawbacks] in [important] messages per split (no blankness or receiver).

Figure 3-4: Synchronous and asynchronous spill ordering

The synchronous spill algorithm ensures that on initial insert $I$ and $a$ adjoined with $a$ at $a$ into PC node are performed in the same order as the corresponding relayed insert $a$ and adjust spill $a$ are performed to the PC with the PC ordering sorting the standard. We can view this requirement around used for the non-PC nodes determines the ordering on initial inserts and adjoind spills, and place the burden on the PC to comply with the ordering.

Suppose that the PC performs initial spill $S$, then removes a relayed insert $a$ from $a$ where $I$, was performed before $a$ at $a$ (see Figure 3-4). We can keep $S$a, compatible with $A_i$ by removing $S$a, overwriting $a$ (while $I$, they'y way in the PC's node). Only $S$a, can be overwritten, not $S$a, can hit revisitin by performing $a_i$ in the PC. Otherwise, $a$'i sky should have been sent to the (thing that a texted. For usually the PC can record an invariate by creating a new record such $a$ kept, and marking $a$ in his sibling. This at this base for the over synchmes for spill algorithm.

**Algorithm.** The entry points are split algorithm at the time and it switches into split algorithm, with the following two phases:

1. When the PC detects that a node needs to serve, it performs the initial split (creates the copies of the core taking etc.) then node released split arrives to the other nodes

2. When a new PC copy receives a relayed split arrives, it performs the relayed split

3. If an PC receives a relayed event and the event is not at the range of the PC, the PC creates an initial event arrives and sends it to the right neighbor

**Theorem 1** *The new synchronous split algorithm satisfies the computer concurrent and related history requirements.*

**Proof** *The split synchronous algorithm can be shown to produce complete and it third balance in the same manner as in the proof of Theorem 1.*

We need to show that all copies of a node have compatible histories. Since relayed events and relayed splits messages are used only consider the same when at least one of the system is an actual action. Suppose that copy $i$ performs a relayed split and PC performs a relayed split. PC has a clean range of PC after $k$ times can be extended for $k$ it. In $k$ as shown earlier, by $k$ relay a copy is in copy form the split history. Suppose that $i$ performs $k$ before $i$ and PC performs $i$ after $k$. If $i$ is on the range of PC, then after a range processing; but the range copy of a copy from the PC a relay a copy; a relay a split. In a range processing in a relay a copy PC is qualified in uphold reading, $i$, for a processing. This is exactly the action the algorithm takes. If

Theorem 2 shows that we can take advantage of the structure of the users and split actors in body mumps replicated states at the utmost codes of the B-tree. In the next section we describe a different type of key-type management which also simplifies implementation and enumerate performance

### 5.2.7 Single-user Mobile Nodes

In this section we briefly recount the problem of lazy node mobility. We assume that there is only a single copy of each node, and thus the nodes of the B-tree can migrate from one server to processor (typically to perform load-balancing). When a node migrates, the host processor can broadcast to all-processors to every other processor that conveys the node (as is done in Exwctarl [36]). However, this type of busy-manner keeps conjectival-neutral effort and doesn't solve the garbage-collection problems.

The algorithm we propose follows the notion immaterial emphasis of the user address. In order to find the neighbors, a node contains links to both its left and right sibling, as well as to its parent and co-children. When a node migrates to a different processor, it leaves behind a forwarding address. If a message arrives for a node that has migrated, the message is routed by the forwarding address. We have left with the problem of garbage-collecting the forwarding addresses (where it is safe to reclaim the space used by a forwarding address). As with the fixed-upon addresses, a child that receives a message for a node that is not-resident can identify the forwarding address, so we could use forwarding addresses[20].

The major algorithm ensures that a forwarding address main unit the processor is guaranteed that no message will arrive for it. Unfortunately, obtaining such a guarantee is complex and implied each-of-the messages passing and synchronization. We want the details of this super algorithm so just goes.

Suppose that a node empties and doesn't have behind it a forwarding address. If a message arrives for the emptied node, then the message clearly has messed tipled. This set-even is similar to the unattended operations as the run-scient. It left protocol, which suggests that we can use a similar message to restore Trap the node. We'll need to find a protocol to follow. If the processor resets a lost node, then that node restores the first leaf to its node to restore former settings. In the over-recovery mechanism is to find a node that is "close" to the destination and follow that set of links.

The other issue to address is the indexing of the columns on the index (once home is only one copy, every node having is constantly compatible). The possible actions are the following: create, split, migrate, and link change. The link change events are a new development so that step we control that in external nodes, and need to be performed in the order issued.

**Algorithm:** Every node has two additional interfaces, a remote monitor and a local. The remote routine allows us to kindly perform related intiative. The local, which indicates the column to a leaf, used to recovery from interruption. An operation is executed by executing an R link low action, as we only need to specify the execution of the return.

*Out of setup:* When a message arrives at a node, the processor first checks if its node is at range. If a check includes routing to tie a Un node level and the message doesn't meet Its Jefd. If the message is out of range or on the setting Jevd, the node causes it at the appropriate direction.

*Migration:* When a node migrates,

1. all actions on the node are blocked until the migration is complete.

2. A duplicate copy of the code is made on a remote processor, [with the exception that the return context assumes by ()]

3. A lock change token is sent to all known neighbors of the node.

4. the original node is deleted

Answer: lattice are performed locally

*Self-safe*  Skill splits are performed locally by placing the sibling to its own prime and assigning the sibling a return position time (to hold that other work is the right neighbor

*Lock charge*  When a code receives a lock change token, it updates the selected link only if the update's return number is greater that the link's current return number. If the update is performed, the new return routine is recorded

*Moving Node*  Its message arrives for a node in a processor, but the processor doesn't store the node, the processor performs the out-of-range action at a locally stored node. If the processor doesn't store a search-assistance node, the action is sent to the root

*Theorem 5*  The bug algorithm can infer the complete, complexity and reduced lattice experiments

*Proof*  There is only a single copy of a node, so its instances are necessarily comparable. Each action takes a point state to a point at every action eventually leads to termination. Therefore, the algorithm produces complete instances

The only unclear actions are the lock change actions. The node at the end of a link can only change that its a split or a migration. In both cases, the node's return

number is incremented. When a link change values occurs at the correct destination, it is performed only if the current number of the new node is larger than the current number of the current node. If the current number is not performed, the node's current number is set equal to that of the new one. Each node starts with a number less than the others, so it ensures to have the first copy of the new one. In this case, it's a link change occurs between two nodes. Let $i$ be at the current number of the current node, $c_i$ the number of the first copy of the new node, and let $n$ be its current number of the current node. If the current number of the new node is $R_i^l = c_i (\frac{n-1}{q_{max}}) + (\frac{q-n}{q_{max}} c_{i,l})$. Then, the history can be rewritten so that it remains valid if.

We state that an implementation of the large single copy algorithm can be executed. By addition in adaptive efficiency and without overhead. The forwarding addresses are not required for correctness, so they can be garbage-collected in a concurrent manner.

### 3.5.3 Scalable Copies

In this version, we assume that leaf-level nodes can migrate, and thus processors can join and leave the computation of the active index. In this case the algorithm we implement a more complex scheme[36]. We assume that the leaf nodes can be replicated, and that the RG of a node even though replication.

The key algorithm that we propose combines elements of the lazy hand copy and migrating node algorithms by using lazy splits, version numbers, and message recovery.

To allow for data balancing, we let the leaf level nodes migrate. The leaf level nodes aren't replicated, so we can migrate them with the lazy algorithm by migrating locks as described in 3.5.2. We count as maintain the RG core replicated in the RG tree property that if a node contains a core leaf node, it must not the set of copies for every node from the

cost to the bed which it does not already help inaccessible. If the processor starts off the last child's node, it organized the set of generators which reinforces the parent (applied recursively). When a processor sizes to require a node refinement, the neighboring nodes are reflected of the corresponding processors with a node change of two. The neighbors link change reduces, we expect to start a node here grows to each a left and right sibling. Therefore, a right-to-point of how a processor starts off the last child's node is always as some other grown further.

We assure that every node has a PE-link cover changes (we can relax the assumption). The primary copy is responsible for performing all actual split actions by registering all grow and assume actions. The join and union actions are ambiguous to the requester either a request or delete. The join and union actions are ambiguous to the requester either a request or delete.

Similarly, a processor wants to assume that it needs some processors to hash its left and right sibling. Therefore, a right request is registering to reconstruct the related action or the node and perform every recovery on all local action requests.

## Algorithm

**Out of range:** If a copy receives an initial action that is out of range, the copy sends the action across the appropriate link. Related actions that are out of range are discarded.

**Insert** 1. When a copy receives an initial insert action, it performs the insert and each copy found across is the other copy, each copy found across is the other copy. Then the copy must hash its left and right sibling it needs.

2. When a non-PC copy receives a related insert, it performs the insert if it is in range, and discards it otherwise.

3. When the PC receives a relayed insert action, it tests to see if the relayed insert action is at large.

    (a) If the card is at large the PC performs the insert. The PC then relays the insert action to all input that passed the replication at a layer earlier that the sensor attached to the relayed update.

    (b) If the card is not at large, the PC sends an natural insert action to the appropriate neighbor.

**Split** 1. When the PC detects that he may a has left, a performs a half split action by creating a new sibling on several processors. Copying most of them to be the PC and transferring half of its keys to the copies of the new sibling. The PC sets the starting session window of one new sibling to be its own session number plus one. Finally, the PC sends of insert action to its source and split the new sibling to be its own right sibling, and relayed split actions to the other copies.

2. When a processor joins a replication of a copy, it sends a join action to the PC that creates the copy. The processor creates the join on that core makes a copy and creates a copy to its neighbors.

**Join** 1. When a PC detects that it contains too few keys, it creates a join action to the PC. The PC detects whenever enough keys at the one that is on each side. The PC relays join actions to all of its copies and sends a join action to the PC of the natural sibling. When the PC receives the join action from the natural sibling it gives up all its keys to the sibling.

**Clone** When a processor clones a replication of a copy, it sends a clone action to the PC that creates the copy. The PC relays a clone action to all the copies of the one that creates the copy. When the PC receives the clone action it makes a clone and creates a clone to its neighbors.

**Expose** When a processor implant a replication of a copy, it sends an expose action to the PC that creates the copy. The PC relays an expose action to all the copies of the natural action. When the PC receives the expose action it makes an expose and relays an expose to its neighbors.

simple action. A mutation the processor flows the list of experts, takes the expert in the other region, and performs a link change action on all of its neighbors.

*Reboot process.* When a core PC copy receives a join on an *update* action, it updates its list of participants and performs the *receive* number.

*Link change.* A link change action is executed using the *targeting-note* algorithm.

*Moving code.* When a processor receives an *actual* action for a node it doesn't man-age, it informs the action to its *'close'* node, or removes the action to the node's *level*.

**Theorem 1.** *The variable region algorithm satisfies the complete, compatible, and or-dered linking measurements.*

*Proof.* We can show that the variable region algorithm preforms complete and ordered linkedm by using the proof of Theorem 2. If we can show that for every node $n$, the history of every copy of $n$ satisfies[1] has a multivalue extension $H$, whose uniform update release on exactly $H_n$, then the proof of Theorem 2 shows that the variable region algorithm also supports the uniform distribution linking measurements.

For a node $n$ with primary copy PC let $h_n$ be the set of update actions performed on PC when the PC has certain number. When copy $c$ is created, the PC replaces its receive number by $c$ and sets it to node's values $h_c = h_n$. Where $H_n$ is the included extension of $h_n$, let $c$ an extract all indices satisfies values in $h_c$ through $c_0$. The PC code releases all the stages of its own replacement number. After a copy $c'$ of number $h_{c'}$ received on update number $h_c$ created at an expand set of the pointer set replacement number. After a copy $c$ is released of a $c$ node all of its updates in $c$. The copy $c$ might perform extra actual update updates consistent with $c$'s primary number[2]. These extra updates apply the actions consistent with the number $c$ expansion consistent between $c$ and its neighbors. Then, the variable region algorithm produces comparable between $h$

Figure 5.1: Interrupter tolerates that its concurrent peers will inject

## 5.3 Conclusions

In this chapter, we have discussed the following:

- Replication Algorithm
- Lazy Updates on a 4G isin
- Consistency Library for Lazy Updates

We present algorithms for implementing lazy updates on a 4G-isin, a distributed E-tree. The algorithm can be used to implement a 4G-tree which never merges empty nodes and performs data balancing on the leaves (we have previously found that the lazy-merge policy provides good space utilization [6] and that leaf level data balancing is effective and fast overhead [26]). We provide a consistency library for lazy updates, so lazy updates to integers can be used to implement lazy updates on other distributed and replicated search structures [14]. Lazy updates, like lazy indices system, protect the efficient maintenance of the replicated index under keyin field synchronization is required lazy updates permit concurrent search and modification of a node, and even concurrent modification of a node. Finally, distributed search structures which use lazy updates are easier to implement than more structures of precision because lazy updates stand the use of synchronization. The next chapter presents the details of our implementation of the distributed E-tree.

# CHAPTER 4

## INFLUENCE FACTOR

### 4.1 Introduction

A distributed environment consists of numerous capable of communicating with each other through messages. We implemented the distributed filter as a general network architecture, a LAN network comprised of SPARC workstations. Every process is capable of communicating with other processes and has sufficient amount of local storage. Each processor acts as a server responding to messages from other processes.

### 4.2 Design Overview

The filter is distributed by partitioning the scalar of the tree across a network of processors. The network of processors represented by nodes in Grid interconnect average moving subtree. To provide a user interface, we integrated Providers in our design. In this design, there is an overall R-tree manager, called the number, which coordinates all the R-tree operations. The number is responsible for routing new processes in different processors who receives incoming *Query* processor is individually responsible for the nodes it is responsible for.

On each processor we have a queue manager and a node manager. The queue manager receives messages from remote processors and maintains them in a queue. The node manager handles message distributes over the neighboring nodes. The distinction of process

Figure 4.1: The Communication Channels

chronization, and the queue manager and the node manager enables bus node manager to be independent of the inter-processor communication method. The queue manager and the node manager is a prominent intermediary for the inter-process communication scheme supported by UNIX, namely message queues (Figure 4.1).

### 4.3.1 Analog Process

The analog is responsible for initializing the Errant. In addition, the analog receives update messages from external applications and sends them to the appropriate processes. Each process is responsible for the decision it makes necessary for the execution it hosts. As the analog implementation, the analog makes the decision it sees at processors and initialized. In order to do so, the analog must have a picture of the global state of the system. The Errant processing will initiate which the analog makes its decision, so the global picture will usually be consistent not at that time. One algorithms take this fact into account.

The worker begins building the tree by selecting a processor (free root processor) to hold the root of the tree. The code manager at the root processor has a socket connection to the worker. Update operations are passed to the root processor and provision down to the leaf level, where the obverse action of the operation is performed. The dashed lines in figure 4.1 represent temporary communication channels established between two processors for the transfer of nodes, which will be described in a later section.

### 4.3.2 Node Structure

Logically adjacent nodes may not reside in the same processor; hence, a parent/child/sibling pointer may refer to a node in some other processor. Also, nodes cannot be uniquely identified by the local address. Every node in the R-tree has a name associated with it that is not dependent on the location of the node. This mechanism for naming nodes is known as location-independent naming of nodes. A typical node would have a parent pointer, the children pointers, and the sibling pointer. In addition to having the highest value in itself, the node must also keep the high value of a logical neighbour. This will enable a node to determine if the operation it cannot be used or instead be either of its siblings.

### Location Independent Naming of Node

Whenever a node is created, it is given a name that is unique among all processors. The attempt, the node name may be a combination of the processor number that creates it and the node identity within the processor. A worker building confining is used to maintain balance across nodes and physical data domain. When a node logically leaves first processor A to processor B, it almost as nest. The advantage of this mode of naming nodes is that a parent-child

so adding node that references the node back onto set from the switch address of back to parameter $B$.

A further advantage of location-independent routing is when the nodes are replicated. All copies of a node in different parameters have the same name. So, the memory and secondary copies of a node can keep track of each other easily.

### 4.2.3. Updates

Our implementation is primarily concerned with the update operations: *create* and *delete*.

#### ▶ Results

As more operations at a leaf processor, assume the key to the appropriate node-key $\alpha$. The insertion of a key causes this node, if the location key $\beta$ that reaches the leaf from the node-key. A corresponding $k$ that causes splits to the original ones, or a deletion which is a split of a data structure after the caused a new set assume is that a data store, by $\alpha$ as the key to preserve the granted is another. The message retrieves the name of the new address $\gamma$ and the modified list and has a deletion of a key $\alpha$ an action-independent leaves the parent operation, external the number of messages in the system, the disk pressure and use, see the $X$ task tree *protocol*.

When this partial node $\beta$ removes a valid message from the child $\gamma$ with a new node-key the node is a also position, with the loop to the action/obligations from the parent processor, instead, we use the $X$ task tree *protocol*.

When a disk node, a create a valid message that the disk child $\gamma$ with the node and a $\alpha$ the node add the legit and a new node-key the child and $\beta$ to the parent processor, instead, we use the $X$ task tree and a $\alpha$ to the parent the node adds the legit and a new node-key the child and $\beta$ to the parent processor, instead, we $\alpha$ to the parent the node adds the child. This key is under taken place over a lot lower level and a with message have the parent $\gamma$ on or on a lot level and a $\alpha$ to the parent the process $\alpha$ child.

means spread till the end. The children of a set are informed immediately of the rule to the parent, so none of the children of a phrase transferred to node $n_q$ will have contents in a second of its up. These identical parent children are updated when a message arrives from the parent up to its child. If the child also fills parent passim as a red resource a message loops up, a red resource node information for this rare up and in the contents is updated its parent passim is the rule to the parent and the resource node's parent up by each, the eldest passes of the parent node $p$.

If a count causes the parent to split, the average prevalence up towards the root node. In the event that the root node splits a new root has to be created. The processes holding the root node causes a new node and makes this the new root. However it informs of either processor that the tree height has increased.

## Delete

The delete operation must treat complications, as delete of a key causes shifting the responsibility of a key range between two nodes. A delete operation requires the key from a leaf node $n$, $z_i \ldots$ node $n$. The terminating action the algorithm takes depends on whether $n$ has implemented a free copy or merge of half $b \ldots$ Merging action represents involves two kinds in terms of synchronization and messages and lines in out and efficient. So, if the neighbors are on the same processor, then the merge is half painless is and allowed as internal is tween simply $z \ldots$ here an empty pointed as needed.

The problem that occurs when pulse can split as well as merge is that error actions can be performed twice at some copies, leading to inconsistency. This transpires when an event occurs at the PC before the split and at a new PC copy after the merge.

When the key ranges of atomic index change due to merging, then care must be taken to synchronise the moves to delete with the splits and merges. Let us consider an instance. Suppose there are three copies of a node at $t_0$, $c_1$, $c_2$ and PC (Figure 4.2). Let the symbol union of key $k'$ be performed at $t_0$, the merged union $c_3$ of copies $c_1$ and $c_2$ at the PC. Before the union $k'$ at PC reaches the PC the index $k$ has split, a key-split, note of unit $a$. The merged union $c_3$ of the PC is forwarded to either both copies $c_1$ and $c_2$. Then the index split of the PC has been.

As the key $k'$ is inserted once and once deleted from all PCs, the union $k'$ arrived at $c_1$ and $c_2$ ahead of the arrival of union of copies $c_1$ and $c_2$ to form the key $k$. But the union of the key $k$ is correct once and once deleted from all PCs in this way. The key $k$ has to be a correct item and once deleted from all PCs, the union $c_3$ had arrived at all takes the merge, then the node $a'$ for which it was intended needs no reverse. An essential feature of a local merge model is illustrated, meaning that all copies receive.

### 1. Free-air scaling

A node $a$ that becomes single does not get noticed until its neighbour updates their index. A processor that receives a sibling empty message blocks deletion and sends an acknowledgement after it has set the link.

Figure 6.2: Stochastic waiters due to merges

After the acknowledgements are received from both employees the spare is fixed. The node patient puts the fix deleted from the parent. A message is sent to the parent node and a marked as deleted. However, the node remains in the fix-fix linked list with its siblings and no acknowledgement message passes the parent. This ensures that no further updates to the node is will be received, or a removed item the fix and the spare is retrieved. In the several before the acknowledgement is received, any operations to the deleted node is are met at the strings (as appropriate). If a node is asked to delete a pointer that does not exist, (or the relayed must be not not served at that copy) but is at its key range, the delete action is delayed until the corresponding count takes mirror. Thus a node has to remember delayed deletes.

## 2 Merge-at-half

In addition to fielding nodes that are empty, it's have incorporated a merge protocol to implement merge-at-half. If the removed of a key reduces it to less than one half its maximum capacity, the node nlmers its keys either to its right or to its left. The idea here is to keep the nodes equally full. If the

right or left neighbor has more than half the boxes, the majority is aligned with the units o...

The transfer of keys between two adjacent boxes, must be recorded in the journal. The parent or male owner of the key range or in child ranlers so that future actions would be checked annually.

When the parent node invokes a message to delete a child node, it removes the pointer to the child. On receiving a change in the key range message from a child, the parent changes the highest value of the child. A change in the parent may cause one of the three situations to take place: a child may merge or sent its own keys to the sibling either merging does not increase. In addition, it all sends the way either or the merge more than one node. If a message is sent to its parent to change the child if the root has only one child of too a root will invoke a message to change the child key to delete and will address sibling at the deletion message....

To initiate a better understanding of text three protocols work, let us look at the algorithms in figure 4.3 through 4.5. When a processor receives a delete message, the message records in the appropriate leaf node and then the procedure 4.3 is invoked. In this algorithm, the key x is deleted from the node x that resides as procedure 4.4. This consists of three steps: if the node of the child procedure by invoking the algorithm Justify_state (figure 4.4). Procedure decide_state may return any of the following values ( 4.4 )

- DIVIDE: In this case the root node has been reached and as the delete process is completed. A Adjust_delete message is sent to all copies of the node.

- **EMPTY-LOCAL**: The parent of node $u$ resides on the same processor, $p$ as node $u$, so the parent is updated if the key siblings at node $u$ and the parent cost cost occurred, (proceed).

- **EMPTY-REMOTE**: The parent of node $u$ does not reside on the same processor, as a message is sent to the processor holding the parent node, indicating that the node $u.cost$ becomes empty and to remove the child parent at $u$.

- **MERGE-KEPT**: The right neighbor of node $u$ resides on the same processor, $p$, so the nodes $u$ and its right neighbor share the keys (merge formation 4.4).

- **MERGE-LEFT**: The left neighbor of the node $u$ resides on the same processor, $p$, so the nodes $u$ and its left neighbor share the keys (merge formation 4.4).

- **MR-MERGE**: If the node $u$ is left to empty size less than half link, then a merge with node $u$ and its right neighbor. If the the right nneigh are updated all the root light and low values of the node $u$. If the parent resides on the same processor. Otherwise, a message is sent to the parent on same other processor to update node $u's$ values. A colon delete message is sent to all copies of the node.

```
Procedure Encryptor(abstract, c) {
Acc = FALSE;
while (!Acc) {
    race = TRUE;
    foo = position of key c in the node c;
    interim[n/2] = c;
    state = Invalidate-shared;
    switch (state) {
        case LISTED;
            conducting deliver (c, c);
            break;
        case RESTARED;
            Enable description deliver (c, c);
            c = original;
            state = VALUE;
            break;
        case MAPPED;
            case REMOTE;
            need Extented mapt(c)!
            state = VALUE;
            break;
        case REMBLENO;
            performinvage.vight(dd)
            break;
        case WARNLIST;
            performavage.cv(dd);
            break;
        case BL_BREAK;
            deliver.old.copy(c);
            need Extented mapt(c) = BL(c)/REQ;
            breakprocessnode(gdc, getwigenow(c), getwigenold(c))
            else
                bindprocessnode(c, getwigenow(c), getwigenold(c), gotwigenold(c));
            containsig deliver(c);
            break;
    }
} }
```

Figure 4.1 Encrypter-Deliver Algorithm

```
Procedure decode _state(E)
 struct node *E ;
 struct node *right  *next
 int count
 if (E->private.priv == IDETAIL)
      return (CONTINUE) ;
 if (empty_node(E))
      if (E->private.priv == COMPLETE)
            return (NEXT  NODE);
      else
            return (NEXT • WORD);
 else {
      if (E->right.priv == IDETAIL && E->right.priv == IDETAIL )
            return (CONTINUE) ;
      if (E->right.priv != PAR_PRIV && E->right.priv != IDETAIL &&
            E->left.priv != PAR_PRIV && E->right.priv != IDETAIL)
                  return (CONTINUE)
      else {
            next = E->next
            right = E->right
            if (next->private.priv != IMAXSIZE -1)
                  state = notAnyplacer
            else {
                  if (notempty (next) • MAXSIZE)
                        next->anyplacer • MAXSIZE)+1
                  else (
                        next->anyplacer • MAXSIZE)+1
                              return (NOUKLIST)  ;
                        }
            else {
                  next = thisList
                  if (prev != NULL){
                        if (notthisright->prev + notAnyplacer != IMAXSIZE+1)
                              state = notAnyplacer
                        else {
                              thisright->prev + NULL
                              if (notAnyplacer+1 != MAXSIZE
                                    notAnyplacer+1 • MAXSIZE)+1
                              if (state • 0)
                                    return (NOOKLIST);  }
                  else {
                        thisright->prev • MAXSIZE
                        thisright->next • notAnyplacer
                        if (state • 0)
                              return (NOOKLIST);  }
            }
      }
}
return (NO_BUG) ;
}
```

Figure 4.1. Procedure Decode State for Entries

```
Procedure preflows_merge_right (c)
shared node *c  ;
oldright = getrightChild(c);
right = getrightChild;
rightType = getrightChild;
parent = merge_right(c);
empty = merge_right(c);
if (empty) {
  /* don't (prune) */
  if (q->query == ITR_NEXT)
    recursive_delete(c, endleaf);
  else
    seed_keyreset.empty(node);
}
node = getrightparent(c);
node.high = getright(c);
if (c->level > LEAF)
  cache=current.childiterate(node);
return.high(c, node);
if (seedleaf.empty(c, node)) {
  if (c->parent.prev == ITR_PREV)
    seed_keyreset.merge(from , getrightparent(c_prev.parent, node));
  else
    seed_keyreset.empty(node);
}
if (node) {
  if (node->separator == ITR_PREV)
    seed_keyreset.merge(from , getrightparent(c_prev.parent, node));
  else
    seed_keyreset.empty(node);
}
}
```

Figure 4.7  Procedure Preflows_merge_right for Deletes

```
Procedure getMerge left merge ()
struct node *p ;
pnr = z->FirstPart;
prevhigh = getHighest(pnr)
snptr = merge.left(  pnr}
if (empty) {
    if (prevhigh.part >= CUR.PART)
       mergetwo.node(prev.part, merge, prevhigh)
    else
       node.rognitem.empty(pnr)  }
else
    if (prevhigh.part >= CUR.PART)
       build(prev.net.curAnd.glptr)
    else
       node.rognitem.rel.gl(pnr)  }
if (z->front < LAST)
    node.n.conait.of(getdata())
sptr = getHighest(snptr) ;
snnd.confirm.merge( )
if (q->getvn.not. >= IVV.PART)
    buildprevsetgetm(n.getglgetall(), getbigh.tell(), getbigh.tell())
}
```

Figure 4.8   Procedure Perform_merge_left for Deletion

## 4.3 Data Balancing the Tree

We have addressed the need for distributing the Stones. Distributing the tree adequately implies that every processor may have many nodes (that is splits). Hence, when there are plenty of nodes in the system, the connection and the rule of binary storage capacity. Hence, the efficient use of storage and other resources, as a necessary to balance the load among processors. We will be discussing the various algorithms for data balancing in this next chapter. However, certain relevant issues, such as methods for dealing with ad of nodes (messages caused by splits) introduced by the underlying network, and low overhead synchronization of non-orderstoring, will be discussed here. Methods for node orderby intended for data balancing and concern sharing are also discussed. We have developed algorithms for dynamic data-load balancing that use the mechanisms of node mobility. In this chapter, we present issues pertaining to the following. Other issues that arise from load balancing are mechanisms for node mobility and use of order information handling.

The fundamental issue in load balancing is the actual process of moving a node between processors. This is termed the node migration mechanism and is common to all of our algorithms for load balancing.

Another important concern is the use of order information that the processors have. Since processors do not have up-to-date information about every other processor in the system, they need rely on the old information to make decisions. In an overloaded processor, either to unload some of its nodes to another processor, or obtain a receiving processor (where the load is explained in the next chapter). As a result, whenever a processor (here for) time will be explained in the next chapter) and follows a regulation protocol to determine the exact the receiver of nodes to transfer. This will be discussed in greater detail in section 4.4.

The node migration algorithm should address the following questions:

1. **Who is involved?** Should the entire and all other personnel be locked up until all persons to the code in transit get updated? If that is no, prohibition would be lost, then should we achieve consistent possibilities?

2. **What is everyone informed?** Does a code but later selected for migration, how and when is every other governor informed of its new actions? Is the system in which each assessment takes place and the other related processes are achieved, what happens in the system that come for the code in transit? How do they get forwarded?

3. **How is Obsolete information handled?** When a code moves it sends an update link message to related processes. Suppose the update message for link change gets delayed and the code comes for a recent issue. The second update message may reach a such a processes before the first one, what happened, should use links to render the problem.

Our objection addresses three problems and possible solutions to those.

In the contextual code visibility, object mobility has been proposed in Emerald [20]. Emerald is an object based language which places emphasis on the mobility of objects. Objects in Emerald can be data object in-process object and the distribution is adaptive to dynamically changing loads. Thus, every object has a forwarding address composed of a timestamp and address. Every time the object moves, the address and the timestamp is updated. If an object moves from node A to node B only node A and node B are updated. When node C addresses the object at node A, the message is forwarded to node B. Finally, node B responds to the message and updates the message back to C with its new address piggybacked. Objects keep forwarding information even after they have moved to another node and are a transient protocol. If no forwarding information is available.

### 4.2.1  Node Migration Algorithm

For the following discussion on the node migration mechanism, let us assume the node messages at a processor wishing to download its nodes has been notified of a recipient processor that is willing to accept nodes. The actual method by which this is done will be explained at a later section.

After the node messages is informed of a recipient for excess nodes, it must decide which nodes to send. This may be based on various criteria of dissimilation. After selecting a node, the node messages begins the transfer. This procedure is explained by providing initiates to the partition panel in the octanehytose.

**Who is involved?** What relates to the first problem is determining an algorithm to be communicated possibly by sending to only the node and the receiver decay whom to accept. In addition to the communication possibly by sending to only the sender and the receiver decay whom to accept. The interaction protocol between these two processors could be done along these nodes to be transferred. There are nodes to be communicated by sending to only the node and the receiver decay whom to accept. The interaction protocol between these two processors could be done along these nodes to be transferred. After a decision has been reached the sender node makes a decision to begin with a point to be transferred. After a decision has been reached the sender node makes a decision to begin with a point to be transferred, so the migration protocol, the receiver initializing a connection and the sender wait and reaches an acknowledgement at a later time.

**What is everyone allowed?** The sender and receiver nodes of no different processors involved in the transferred nodes. The sender node then decides that it has been sent it again or it timed out an acknowledgement as performed or performed at this point of the transfer. If no acknowledgement as performed at this point of the transfer or performed, the receiver must have his own receiver. The receiver can them begin to the operations are performed at their nodes in the vicinity of a sender.

**When is everyone coherent?** The sender and receiver messages arrive for the code at the same time and the sender first received the code at the same time and the receiver first received the message arrives for the code at the same time and the sender first received.

| Processor 1 | Processor 2 | Processor 3 | Processor 4 |
|---|---|---|---|



local front
of Processors

code, a record
from A to B

Message 1 sent
from B to F

code, a vector again from
B to C

Message 2 sent
from C to D

Figure 4.1 *Static Mappings*

## 4.4 Negotiation Protocol

In all our algorithms, no extra resources are sent to reduce the slope parameters of a charge at the current values of a processor. Thus, if the number of nodes in a processor increases or decreases due to splits or merges, other processors are not aware of it. Nodes or the nodes' parents informed about those changes, the cause being that a merge may be able to merge two nodes into a single node and an effective solution to all our problems. Thus, an the local informations about the slope values and we want to set up our algorithm so that it does not dominate.

In this algorithm, we choose heuristics practiced for the negotiation. During the actions that the nodes' processor decides to share some of its data and a receiver processor is chosen, either by the nodes we considered or by need to inform the receiver processor about those values it is able to share as the data for the receiver. If the node is split or merges, the set of all information becomes of the negotiation protocol.

We have designed an choose heuristics practiced for the negotiation. During the actions that the nodes' processor decides to share some of its data and a receiver processor is chosen, either by the nodes we considered or the consider processors. Once this algorithm is completed, node transfer takes place. It should be noted that an messages are sent to other processors informing them of the migrations or change in the count of the nodes and the count of the processors.

Finally, we have ported our implementation to the SDK a shared memory and threescore machine with 80 processors that supports message passing by providing

SKID markers [DX]. The portray of our implementation shows that our systems is portable and easily scalable to a large number of processors.

## 4.6 Conclusion

To conclude, this chapter has addressed the following:

● Design reuse for implementation

● Data balancing the dll-tree and the fundamental primitude necessary

● Portability of the implementation

In this chapter we have discussed the implementation of the distributed ll-tree on a network of Sparc stations and the processes needed to manage the dll-tree. Update operations, insert, search and delete are performed on the ll-tree. We have presented how these operations are performed and what complicatedes the delete operation present and how we overcome there. To facilitate data balancing on the distributed ll-tree, we have introduced the conversion of memory codes so that a code always lie across between processors. We will see that this mode varying alse a useful when stabilizing codes at various processors, cases of copies of a code have the same name.

We have presented two mechanisms fundamental to data balancing, namely the node migration algorithm for the actual movement of nodes between processors and the equilibrium protocol to overcome the effect of contented archivations. Methods by which our algorithms and protocols tolerate rate of rarely overcome embedded because of intrinsic delays are also presented.

Finally, to study the portability and scalability of our implementation, we ported it to the KSR, a large tools shared memory multiprocessor system. In the next chapter we resume the algorithms for reduction and load balancing and present performance results.

# CHAPTER 5
# PROVENANCE

## 5.1 Introduction

In this chapter we present the various algorithms for replication and data balancing and discuss their performance in detail. Experiments using the two examples for replication, namely *full* replication and *path* replication were conducted. Results show that *path* replication will result in suitable distributed B-tree. We tabulated the tree scalability by simulating a large scale distributed B-tree and performing largescale experiments on it. Several load balancing algorithms on the AD tree. These algorithms from, namely random, *merge* and *appraise* merge algorithms, have been developed for data balancing on the AD tree, and three or four that appraises algorithm makes the B-tree scalable. The experiments have been conducted on single computer and the results of running of B-tree. We simulated the tree scalability by simulating a large scale distributed B-tree and perform largescale experiments on it. Several load balancing algorithms can be best exploited on large machines. The balancing algorithm of the AD tree offers the best and The first load algorithms are scalable in large numbers of processors

We conclude this chapter by comparing the performance of the different load balancing algorithms using the different data structures. We present no different results and the best load balancing algorithms can best be exploited on large numbers of processors. We find that both algorithms are scalable in large numbers of processors

## 5.2 Replication

In this section, we describe two algorithms for maintaining consistency among the copies of nodes. Based on the theoretical framework presented in Chapter 3, we

91

have incorporated this replication strategies in our implementations. Our implementations of the Fixed Partition copies algorithm is termed **Fxd-Replication** and that of Variable copies is **Prtl-Replication**. We will briefly discuss the algorithms and the implementational issues in sections 5.3.1 and 5.3.2.

When the value of the $K$ sets are replicated, an obvious concern is the consistency and coherency of the various replicated copies of a node. Subsections 5.3.3 will present the mechanism by which our implementations maintain coherent replicas.

### 5.3.1 Fxd-Replication Algorithm

The Fixed partition copies algorithm (FPC) assumes every node has a fixed set of copies. An most operation searches for a leaf node and performs the insert action if the leaf becomes full, a leaf split takes place. In this algorithm, the Primary Copy, PC maintains all leaf nodes and sends a delayed split to the other copies. Any initial starts of a root PC copy or a split overflow happens and at each node through the delayed split.

In our implementations, the $K$ sets is distributed by having the leaf level nodes at different processors. Leaf level nodes are not replicated and only those nodes are allowed to migrate between processors. Whenever a leaf node migrates to a new processor (not thus currently shares no home), the entire leaf of its non are replicated at the processor. Consistency among the replicated nodes is maintained by the primary copy of a node sending messages to all the copies.

Since the entire tree has been replicated at the beginning of the load level nodes at different processors. Leaf level nodes are replicated and only those nodes are allowed to migrate between processors. Whenever a leaf node migrates to a new processor and after tiles of the same resources are replicated at the processor. Consistency among the replicated nodes is maintained by the primary copy of a node sending messages to all the copies.

- **Algorithm:** The Fxd-Replication replicate the tree at nodes when a processor (reader) downloads some of its leaf level nodes to another processor (receiver). After the leaves are transferred, the reader checks to see if its current has entered

### 5.5.2 Path Exploration Algorithm

- **Algorithm.** One operation for path replication is synchronous, based on a handshaking protocol. When two processors have estimated in the load balancing protocol, a decision has to be made concerning the path from the root to the captured leaves. Either the sending or receiving processor can ensure that the path is sent to the receiver. In our algorithm, the receiver determines what another nodes are needed after removing one leaves. It then sends requests to the processor holding the primary copies of the antitons to get the paths. As the receiving processor takes the responsibility of obtaining the path, the sending processor is free to continue. The receiving processor cannot do usual copying work in modern chip build on bottles to make it's active or receiver. Once the path is obtained, the receiving processor can handle operations (inserts and searches) on its own.

### 5.5.3 Replica Consistency

The operations like content replication/makes handles are number and counts. A search operation is the same as in leaves operation, except a key is not checked. A search returns a version in before as there are issues and further related message to be issued. An operation in the distributed B-tree can be cantained on any processor. Since the orders have only to be at multiple of operations on it, changes in a node copy in any processor must be a node copied even as a split, or all operation must be obtained to make the split. This is done in the following way.

- Insert: An insert operation can be performed on one copy of a node. After inserting the entry, the processor sends a **enlarged insert** to all other processors

that hold a copy of the cache. When a processor removes a relayed insert, it performs the reset operation locally.

- Rule 4 *apply operation* is first performed on a leaf. If the local parent exists on the same processor, the *split* is achieved at the local parent. If the parent does not exist locally, then a **relayed split** is sent to all processors that hold a copy of the parent node. Otherwise, a **relayed insert** is sent.

## 5.5.4 Preliminaries

Here, we compare the performance of full replication and path replication strategies for replicating the *w-fan* nodes of a B-tree.

**Distribution, Speeds and Operations   Experiment Description:** In the experiment, 10 000 keys were inserted; insertion rates patterned as 5000 key intervals. The B-tree is distributed over 4 to 16 processors. Each node in the B-tree has a maximum fanout of 5, and average fanout of 3. We observed the number of times a *split* request has been made by a processor, the number of times that a *leaf last-insert* request has to be messaged (in most overloads), with priority being given to the *path* request. We collected statistics as to how many *consistency messages* are needed to maintain the distributed, replicated B-tree were being measured. In order to make this estimate more prominent, and finally how many *index* messages are needed to search at the end of the run.

The *Message Overhead* (Figure 5.1) graph shows the number of messages needed to maintain the replicated B-tree. We see that in case of full replication, the number of messages for a 4 processor B-tree is around 7600 and for 16 processors it is around

98



Figure 5.1. Full versus Path Replication: Message Overhead

35000 (i.e. the message overhead has increased linearly as the number of printers sent). However, for a path replicated 3 lists, for 4 processors around 3000 messages are needed and for 12 processors only 8000 messages are needed: not even a linear increase.



Figure 5.2. Full versus Path Replication: Space Overhead

The Space Overhead (Figure 5.2) graph shows the number of nodes stored at all processors at the end of a trial. The graph is similar in nature to the message overhead graph. In this graph we consider only the index nodes that account for the vast

storage at each processor, the leaf nodes remaining nearly the same for all processors) as the number of processors increases. The full replication, we see that for a 4 processor B tree the number of active nodes stored at 1302, whereas for 32 processors B tree the number of nodes is 5202, a nearly three-fold increase. It can at a quick explanation B size, this number of active nodes stored over the entire tree for 4 processors is 980 and for 32 processors is 1206, not even a two-fold increase.



Figure 5.6 Patch Replication: Width of Replication at Level 2

The Width of replication at level 3 (Figure 5.2) graphs those how nearly level 3 index nodes are replicated at each processor for a code replication B tree. We attained level 3 since activity takes place at the leaf level 4, and affects mostly at level 3. This bar chart shows the number of nodes in the B tree that have a square, where a more from 1 to 5, with the concentrations being under with 1 copy at 9 processors. The other chart, explicit of replicated nodes over all processors, shows that, even as we increase the number of processors, the level 3 order nodes are not widely replicated at all processors, with there being 981 copies for a 4-processor system and only 941 copies for 32 processors.

**Fault-tolerance** causes her restructuring overhead, her can require a recent ms wait many processors for its transition. We measured the number of hops required for the search phase of the recent processor after 5000 inserts were requested on an 5-processor search, and path explorers required 1.23 messages per search (additional overhead of 12 messages).

Does the above observations, as we see that a path exploited distributed 50-ers performs better than a fully exploited one not a highly scalable (Figure 5.3).

## 5.3 Data Balancing

We have performed data balancing on the 50-ers and the 50-ers. We will discuss the algorithms and the performance of the two separately.

### 5.3.1 The 50-ers

The results obtained from the implementation of a replicated 50-ers led us to explore other algorithms for data balancing on a replicated 50-ers. The experiments with the replicated algorithms led us to conclude that a path exploited 50-ers was more scalable than a fully replicated 50-ers. Hence, we simulated a path-explorated distributed 50-ers. Our objective is to develop data balancing algorithms and also to observe their performance and overhead incurred.

#### Algorithm

In the current design, a fault is placed on the maximum number of nodes of the tree that a processor can hold, termed as the **threshold**. In addition each node has a split level ($0, 1$ threshold) as the number of nodes. This represents a warning level indicating a need for distribution of the nodes. Whenever a node splits, the current number of nodes it should append to and then inherit. If the current number of

nodes exceeds the soft limit, the processor must distribute some of the nodes it has to some other processor. Our algorithms are characterised by the method by which the receiver processor is selected.

- **Centralised Data Balancing**

  One approach is a centralaralised one, where the analyst is responsible for choosing the receiver. The centraliser processor requests the analyser for another processor that can share its current load.

  The earlier but intalisied information about all processors' current capacity. Based on that abstrict information, the analyst selects a receiver processor.

- **Distributed Data Balancer**

  Another approach would be to leave the decision making to the individual processors i.e. a distributed data balancer. Here, a processor wanting to download some of its nodes picks other processors for load information. The picking can be done in two ways. We assume that every processor has a list of participating processors in an array.

  - **Sequential Picking** : Here, processors begin probing other processors in sequentially in their own load. A processor will pick the first processor after itself and checks if that processor has sufficient capacity. If not, it probes the next processor in line, and so on.

  - **Random Picking** : In this approach, the processor (randomly probes other processors. The randomly selected processor is checked for available capac-ity. If it does not have enough capacity, then another processor (randomly the previously rejected ones) is selected randomly.

After a written agreement a law firm selected, the smaller a real estate owner a —— selected by a negotiation proposal [ ? ]. In this proposal, they decide exactly how many tasks are to be transferred from the worker a to the minister c. The negotiation proceeds in a very rational manner, in the interests that the smaller a receives a real estate and the actual tasks transferred and the actual tasks transferred to the minister c reaches may represent over again and issue a change in their interests. Also, in the case of the real estate load balancing proposal, once the worker has out of data information about each processor's status the regulation works well because of the regulation proposal.

## Preliminaries

The performance of the δβ-axis and the δα-axis depends on how the nodes are distributed among the processors which is not depends on the δ-axis balancing of parallel. In addition, data balancing occurs in one overhead.

There are many non-algorithmic factors that can affect performance. First, the number of bugs that an operator requests to find in data increases with the length of the iterator. Secondly, it really is unfortunate when there is an unbalanced load among the processors over the data δβ-δα. Finally, the intrusion occurs which additional change is made available in the search structure affects the performance of the δα-data balancing algorithm. To reduce the number of parameters we need to separate out expressions and the following two scenarios:

• Incremental Growth.

When the change to the distributed nodes can low, the system manager need add storage capacity to some of the processors, or close the δβ size to spread to more processors. Presumably, an ordinary incremental storage growth of the processors that close the δβ size. If an equivalent in adding a disk to a node or running a new storage site. When a processor wishes to store some of its

weeks, and all the currently active protons so, may thus threshold, a new precious may be tested up, or in the event that the precursor have a marked, a measure is selected suddenly and threshold is summed by a function of its current capacity. The overlasted precursor then short is rate with the new precursor with newly added capacity.

• Fixed Single Data Scheme-sp

To study the effect of large losses on the width of exploration, we fixed the length of the year for all of the increments.

To determine the nature of a large scale effects, we made a simulation study of these fractions on a dB tier. We compiled the number of increase logs reached in complete an operation, and the width of exploration, in average number of replies of a next. We are overly convinced with the width of exploration of least 8 nodes (which are most of the index nodes). The width of exploration is a measure of the mean overhead of maintaining a distributed index.

**Experiment Description.** We count **Experiments, Results and Discussion.**  as noted Notes with a uniform random distribution of keys. After the initial 8 keys is created we vary the key distribution against dynamically. To study the effect of the load balancing algorithm thus the distribution changes we distribute the key to a new precursor position, where we enumerate the keys to a new precursor, thereby losing about 16% of the requests to re-entered in one of our test processors.

To study the load surprise balance order execution, we collected distributed amounts at increasing degree, which from 12000 keys inserted in the S tier to each transaction. We count the processors capacity access of the number of keys.

storage width of replication = # of copies of interior nodes / # of interior nodes

A finer violation is the width of replication at each level.

storage width of replication at level $i$ = # of copies of level $i$ nodes / # of level $i$ nodes

We found that the width of replication is affected significantly by the clustering which leaf nodes in majority. We first used random selection, where a generator that has to decimate in load clusters (in leaf nodes emulation). With this we found that the number of replicated copies was large. So, we improved upon this by rooting out all leaves of a parent node. That is, we selected leaves sequentially. The results obtained were much better and we present them below.

## Results

- Centralized and Distributed Data Gathering



Figure 5.4: Performance of Load Balancing

The *Performance test* charts (Figure 5.4) show the polynomial capacity after the insertion of 140 000 keys. Where the "hot spots" distribution is used with node

**Table 5.1 Load Indexing Statistics**

| Event | Average Number of Buses | | | Average Number of Buses | | |
|---|---|---|---|---|---|---|
| | Centralised Load Index | Distributed Load Index | Uncontrolled | Centralised Load Index | Distributed Load Index | Uncontrolled |
| | | | | | | |

In Section 5, processes 3 and 4 are the load preemptors, and contain a disproportionate number of records. Without load indexing, the processor would incur a serious performance load (using array buffer file filter). One load balancing algorithm distributes the records evenly across the processors.

This is where the calculation of load indexing is required. The calculation of the overhead is such that the records are distributed across the processors. The event could be executed in parallel using the processors.

Table 5.1 shows the calculated average number of buses required to load each of the processors. All uncontrolled rows over 1000 buses (about 40% of the maximum).

Table 5.1 shows the calculated average number of buses to load each of the processors, with the distributed load index being the largest over 1000 buses over the centralised load index (about a 20% reduction in performance).

Note that under the uncontrolled scenario a larger number of buses is required to load the processor. With over 1000 buses required to load the whole system, a larger number of processors is required to achieve the same performance between load index and uncontrolled. Using the uncontrolled scenario reduces the behaviour of a processor when the count of the processor is only a single processor. For example, considering one processor under the uncontrolled scenario, the capacities of the processors.

the essential probing and the random pricing mechanisms, the random position-ing seems to impart a little less number of probes than the sequential one. In sequential probing, a processor (say X) may be probed by four neighbors from the two (say I and J), but can only serve one. Hence, the actual processor (Y) may have to probe another processor before its request is met.



Figure 7.3: Average Number of Hops/Search

The graph of the average number of hops/search in processors (Figure 7.3) shows that even as the number of processors increase, the hops/search stays more or less constant. The number varying from 1.8 for sequential to 1.7 for random. In the case of sequential, this corresponds with the distribution of the index number.

We also plotted the average number of hops/search when the number of searches was varied. It is seen that the more the searches the lesser is found and it is more often the case. As more searches increase, the number of hops decreases.

Figure 3.6 Width of Explanations at Level 2



Figure 3.7 Width Of Explanation

The width of explanation charts (Figure 3.6, 3.7) show that in the case of no load-balancing, the width of explanation is lowest, while all the load balancing algorithms move about the same width of explanation. We have examined the width of explanation at level 2 since I saw activity takes place at this level into above the shell. All the above graphs show that the width of explanation is small, on average of about 2. Sequential selection of leaves has a lower width of explanation than random selection as previously mentioned. The width of

replication over all levels a 4-9 by 28 processors and with node average fraction of 7 when leaves are selected randomly, while it is only 2.9 when leaves are selected sequentially. Similarly at level 5 at 3.1 for random selection and 1.7 for sequential selection of leaves

All the above results indicate our load balancing algorithms perform very well in maintaining a good and very close tree balance among processors. The distributed random pairing algorithm requires lower number of probes and moves than our other algorithms. The distributed algorithm with sequential probing reduces the number of hops per search and width of replication tree in close the others. Also, the sublinear increase indicates that this algorithm is suitable for scaling to large trees with large leaves and tree entry processors.

• Incremental Growth Rate Balancing

As explained above in this algorithm, when some of the processors have available capacity, instead of increasing the capacity of every processor we select a processor randomly and activate its capacity. The results obtained also similar pattern to that of the general algorithms.

The graphs in figure 4.6 show that the average number of hops varies between 1.9 and 4.6 as we increase the number of processors from 16 to 28. The width of replication in level 5 (Figure 4.6) is about 1.7 and the width of replication over all levels (Figure 4.6) varies from 5 to 2.7

• Fixed Sample Tree

We performed simulation on fixed sample tree B tree, by running up to 2.8 million hops and varying the average branch from 10 to 60 (average branch is 40% of the maximum branch [8]). In the first experiment we fixed the tree height

Figure 5.9  Incremental Growth Algorithm: Average Number of Steps/Search



Figure 5.10  Incremental Growth Algorithm: Width of Replications at Level 1

Figure 3.18: Incremental Genetic Algorithm: Width of Exploration

in $d$. When the cost of the tree had this desired average factor we collected statistics. We noted the processors capacity or extent of the number of loxere is fine, the number of solve level nodes, and the number of locis. We also noted the number of times a processor creates the load balancing algorithm, the number of probes required, the number of nodes that it transfers and the average number of times a load scale leaves between successive halter with respect to the nodes at the solver level. The primary of these statistics has been distined as the center of small locust trees [20], so here we concentrate on other important statistics such as the average number of average loces for a search and the average width of exploration at level 2.

— Fixed height of a tree:

the first performed experiments with fixed height trees of 4. The graphs in the figures 3.11, 3.12, 3.13, 3.14, 3.15 and 3.16 show the width of explorations at level 2 and the width of explorations over all levels, plotted against an increasing branco for a fixed number of processors. The graphs

(iii) Enough A.13 show the variation of the number of keys with factors for a fixed number of governors

Figure 5.11 Depth $\times$ Tiers: Width of Replication at Level 2 for 18 Processors



Figure 5.12 Depth $\times$ Tiers: Width of Replication at Level 2 for 18 Processors

Figure 5.15  Height of Tree  Width of Exploitation in Level 2 for 16 Processors



Figure 5.16  Height of Tree  Width of Exploitation for 32 Processors

Figure 5.15 (Right 4 Tier) Width of Replication for 16 Processors



Figure 5.16 (Right 4 Tier) Width of Replication for 32 Processors

Figure 5.17. Height 4 Tree: Average Number of Bags/Search for 16 Processors



Figure 5.18. Height 4 Tree: Average Number of Bags/Search for 32 Processors

Figure 3.18  Weight 4 Tree: Average Number of Pings/Search for 38 Processors



Figure 3.19  Weight 4 Tree: Variation of Average Number of Maps/Search with Processors

Panel 6?



Figure 5.11: Height 4 Time: Variation of Width of Exploration at time of production
ments

Panel 6? Protozoa II



Figure 5.10: Height 4 Time: Variation of the Width of Exploration of by level.

Figure 3.19 : Single a Tree Linear Regression of the Width of Replacement...

The MOB at level 2 reaches a plateau around 2.1 for 16 processors (2.11) around 2.6 for 20 processors (Figure 3.12) and 3.2 for 18 processors (Figure 3.17). Similarly the width of vegetation over all levels reaches their last 16 processors (for plateau is 3.4, 0.79 (Figure 3.14), for 16 processors it is 3.1 (Figure 3.16) and for 20 processors it is 2.9 (Figure 3.17). Again, the number of busy quality reaches a plateau. From the table 3.2 we see that the number of busy a multi instinct with increasing Variant, and reaches a value of 1.95 for 16 processors

For a comparative study of the graphs in figures 3.11 through 3.18, we have combined the data into a table 3.3.

From the table 3.3, we observe that for a 1B tree with a large Variant, the width of vegetation and the number of large set openings depend for the number of processors only. Therefore we can predict the number of busy and the width of vegetation by exploiting the pattern in the chosen value only in increasing number of processors.

• Number of busy

Since the table 3.3 and figure 3.18 we see the effect of increasing the processors on the number of busy. Our results indicate that the hops do not increase significantly well and reach only a value of 1 at 30 concludes that in a large scale diffuse with 1 levels, an average busy

Table 17: data for Excel length of t different

| Statistics | Processors | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| FULL1 | 10 | | | | | | | | | |
| | 20 | | | | | | | | | |
| | 30 | | | | | | | | | |
| | 40 | | | | | | | | | |
| | 50 | | | | | | | | | |
| SQRT | 10 | | | | | | | | | |
| | 20 | | | | | | | | | |
| | 30 | | | | | | | | | |
| | 40 | | | | | | | | | |
| | 50 | | | | | | | | | |
| Non | 10 | | | | | | | | | |
| | 20 | | | | | | | | | |
| | 30 | | | | | | | | | |
| | 40 | | | | | | | | | |
| | 50 | | | | | | | | | |

of Metcal Barchand over 20 processors, at most 2 hops per operation are required.

• Work of Replacemen

In figure 3.14, we plot the glorious value of the width of replication at level I against the number of processors. The linear regresson of the data shows that the slope for the uniform pairing data at 1000 for the sequential sorting algorithm, the slope at 5001. Based on the formula, with $k$ operations at level $I = 1000 = 1000$ we recalculated the width of replication well as figure 1.32 we show a comparison of the two set of values. We see that the experimental values are nearly then obtained theoretically. So, it has been a 1000 processors and a biases of 1000, then the ROM bit level I nodes as about 22 for uniform pairing and 32 for sequential sorting.

Another interesting characteristic is the radiance of the width of replica-tion seak the level of the test. At the peak explosive-depending for the $\Delta$2-cises, the width of replication for the next is the number of processes, next for the domain is 1. We present the WOR number a plateau seek moraning domain for a fixed number of processes and the number of legs is nearly constant, for a fixed number of processes. However, the WOR, at level 2 is higher, reaching a value of $8.71$ for $32$ processes. The WOR over all levels is $7.15$ for $32$ processes. For the length, the WOR at level 2 is at level $3$ and the level $2$ is lower than the level $3$ as a height loss. The width of replication at this that half for the fixed level.

Fixed height of $2$-trees

For fixed height $2$-trees, from the double $2.21$ through $4.31$, we notice patterns similar to that of height $1$-trees. We notice that the WOR number a plateau seek moraning domain for a fixed number of processes and the number of legs is nearly constant, for a fixed number of processes. However, the WOR, at level 2 is higher, reaching a value of $8.71$ for $32$ processes. The WOR over all levels is $7.15$ for $32$ processes. For the length, the WOR at level 2 is at level $3$ and the level $2$ is lower than the level $3$ as a height loss. The width of replication at this that half for the fixed level.

Fixed height of $2$-trees is constant, it decreases lower the WOR at level $2$ goes level lower than this the realty to $9.25$ and $3.71$ two level constant higher than this two-levels the $9.25$, releases lower the WOR at level $2$ goes level lower. For fixed height $2$-trees it decreases lower the WOR at level $2$ goes level lower than this the realty to $9.25$ and $3.71$ two level constant higher than this two-levels the $9.25$.

On the other hand, as the number of processes and obtained a domain for the loads of radiance at level $2$ is $7.09 = 10.5 \times P$. We understand the results of replication at level 2 seap the radiance of legs at level. A level $3$ at this point for these number of legs are nearly the same. For more an experiment and theoretical values obtained. Again, we conclude that the WOR and the number of legs are greatly affected by the number of processes over which the $8$-tree is distributed.

Figure 5.24 Height it Term Width of Implantation at Level 2 for 15 Processers



Figure 5.25 Depth it Term Width of Syrflexion at Level 2 for 30 Processers

Figure 5.26 Height 2 Year Width of Replication at Level 2 for 3I Processes¹



Figure 5.27 Height 2 Year Width of Replication for 3i Processes

PLATE # 34



Figure 1.37  Deple 1 Tree  Width of Explosion for 3d Premises[?]

PLATE # 35



Figure 1.38  Deple 1 Tree  Width of Explosion for 3d Premises

Figure 5.20. Height 2 Tree: Average Number of Bugs/Search for SB Preservery



Figure 5.21. Height 2 Tree: Average Number of Bugs/Search for SI Preserves

Figure 3.21: Height 1 Tree: Average Number of Keys/Search for 30 Processors



Figure 3.22: Height 1 Tree: Linear Regression of the Width of Replication

Fixed height of V trees.

For fixed height trees of height $h$, we could not get statistics for factories larger than 30, as for tree was very big. The WGR at level 2 is 2.33 for 30 processors and the WGR over all levels is 7.64 for 30 processors. The number of hops is enormous at 2.33. We include the charts: 3.36 through 3.42 and table A.6 for the sake of completeness.

Table 5.3 data for Fixed height of 0.49 mm.

| Iteration | Thickness | Stress | | | | | |
|---|---|---|---|---|---|---|---|
| | | 10 | 15 | 20 | 30 | 35 | 40 |
| Phase 1 | 10 | 3.73 | 3.55 | 3.73 | 3.58 | 3.14 | 4.47 |
| | 20 | 3.90 | 3.73 | 3.55 | 3.73 | 3.18 | 4.48 |
| | 30 | 3.98 | 3.84 | 3.90 | 3.85 | 3.58 | 3.72 |
| | 40 | 4.89 | 3.80 | 3.90 | 3.85 | 3.55 | 3.73 |
| | 50 | 3.85 | 3.90 | 3.85 | 3.58 | 3.73 | 3.55 |
| | 60 | 3.58 | 3.73 | 3.85 | 3.90 | 3.72 | 3.58 |
| MTR | 10 | 3.85 | 3.58 | 3.73 | 3.55 | 3.73 | 3.85 |
| | 20 | 3.72 | 3.85 | 3.58 | 3.73 | 3.55 | 3.73 |
| | 30 | 4.14 | 3.90 | 3.85 | 3.58 | 4.45 | 3.72 |
| | 40 | 3.90 | 3.85 | 3.58 | 3.73 | 3.55 | 3.73 |
| | 50 | 3.58 | 3.73 | 3.55 | 3.85 | 3.90 | 3.85 |
| | 60 | 3.73 | 3.58 | 3.85 | 3.90 | 3.72 | 3.58 |
| A-por | 10 | 3.85 | 3.90 | 3.58 | 3.73 | 3.55 | 3.73 |
| | 20 | 3.72 | 3.58 | 3.85 | 3.90 | 3.85 | 3.58 |
| | 30 | 3.73 | 3.85 | 3.90 | 3.58 | 3.73 | 3.55 |
| | 40 | 4.14 | 3.73 | 3.55 | 3.85 | 3.90 | 3.72 |
| | 50 | 3.58 | 3.90 | 3.85 | 3.73 | 3.58 | 3.73 |
| | 60 | 3.85 | 3.58 | 3.73 | 3.90 | 3.55 | 3.85 |

Figure 2.32: Graph 1: Error Width of Implication at Level 3 for 10 Processors

Figure 5.19. Graph 1 Tree: Width of Explosion at Level 1 for 30 Provider(s)



Figure 5.20. Graph 1 Tree: Width of Explosion at Level 2 for 30 Provider(s)

Figure 5.31  Length × Time: Width of Rectification for 10 Pressure(s)



Figure 5.32  Length × Time: Width of Rectification for 20 Pressures

Figure 5.38. (Depth-3 Tree) Width of Replicomm for 30 Processors



Figure 5.40. (Depth-3 Tree) Average Number of Trees/Branch for 30 Processors

Figure 4.11: Height 5 Tree: Average Number of Hops/Search for 50 Processors.



Figure 4.12: Height 6 Tree: Average Number of Hops/Search for 50 Processors.

Table 2.x. Data for Fixed height of 2.40 cm.

| Variable | Precision | Errors | | | |
|---|---|---|---|---|---|
| | | I | II | III | IV |
| HOR-X | 10 | 1.20 | 1.74 | 3.45 | 4.68 |
| | 20 | 1.61 | 2.34 | 3.90 | 4.93 |
| | 30 | 1.95 | 2.60 | 3.97 | 4.97 |
| | 40 | 2.42 | 3.06 | 4.07 | 5.00 |
| | 50 | 2.86 | 3.49 | 4.13 | 5.04 |
| | 90 | 4.37 | 4.80 | 4.96 | 5.25 |
| VER | 10 | 1.20 | 1.74 | 3.45 | 4.68 |
| | 20 | 1.64 | 2.38 | 3.94 | 4.96 |
| | 30 | 1.95 | 2.60 | 3.97 | 4.97 |
| | 40 | 2.42 | 3.06 | 4.07 | 5.00 |
| | 50 | 2.77 | 3.40 | 4.09 | 5.03 |
| | 90 | 4.37 | 4.80 | 4.96 | 5.25 |
| Mean | 10 | 1.20 | 1.74 | 3.45 | 4.68 |
| | 20 | 1.62 | 2.36 | 3.92 | 4.95 |
| | 30 | 1.95 | 2.60 | 3.97 | 4.97 |
| | 40 | 2.42 | 3.06 | 4.07 | 5.00 |
| | 50 | 2.81 | 3.44 | 4.11 | 5.04 |
| | 90 | 4.37 | 4.80 | 4.96 | 5.25 |

Table 3.3. Comparison of Fixed height 3, 4 and 5 over with Errors 20 and over 50 processes.

| Height | Width at level 3 | Width | Number of items |
|---|---|---|---|
| 3 | 2.40 | 4.08 | 1.89 |
| 4 | 2.52 | 4.40 | 1.60 |
| 5 | 1.97 | 4.04 | 2.15 |

135

In our code we observe that the quality obtained bears the variability of the distributed R-tree. The number of keys depends on the length of the tree and the NGM depends as the number of processors, but grows very slowly. From the comparison table of the different lengths at table 3.5 we can see how the width of replication and average number of keys per operation vary each too.

## 3.6.3   The dir tree

The dB tree is a practical distributed index constructed from the distributed B-tree. The main purpose of this being to reduce the communication cost by storing fewer leaves and thus saving less overhead. We have observed that, instead of maintaining separate leaves for the connection keys stored in a processor, a more efficient approach is to maintain a single leaf (i.e. an extent) that maintains key range information only, and then the keys in a local data structure.

The difference between the dB tree and the R-tree is that it is the leaf balance that decides whether to split or merge a leaf node. The leaf balance is invoked when a key is moved into a leaf node. If the leaf balance describes that the previous leaf has too many keys, it decides to download some of its keys to some other processor. It selects a leaf node and decides to either perform a merge or a split. The processes each select to merge or give away the info sliding to also involved based on certain criteria. We redirection for the whenever routine for the leaves and the processors before.

## Algorithms

In each of these algorithms, the leaf balance decides if the processor has an unbalanced set of keys. Let the mean number of keys be k.

**Synthesis.**   As the mean suggests the leaf node to be merged or split is selected randomly.

- Step 1. Push a suction valve. If the node $n$ is o-a-c-k by processor $P$ then go to step 2, otherwise go to step 3.

- Step 2. If the node $n$ has a right neighbor $n$ and it's cursor processor has available capacity, then transfer the cursor into this node and stop.

- Step 3. If the node $n$ has a left neighbor and it can be merged with relieve its left-m-right neighbor. If there is no such left node, then the largest extent that is chosen for a unit.

  - Step 1. Scan through the list of nodes until an extent is found that is owned by the processor $P$ that lies at least $K$ keys. Let this node be say $n$.

  - Step 2. If the node $n$ has a right neighbor $n$ and it's cursor processor has available capacity, then transfer the cursor into this node and stop.

  - Step 3. If the node $n$ has neighbor left has available capacity, then transfer the cursor into this node and stop.

  - Step 4. If the node $n$ has a processor that can take all the excess keys $K$, then move the cursor into this node and stop.

- Step 1. Push a suction valve. If the under-o-a-c-k by processor $P$ then go to step 2, otherwise go to step 3.

- Step 2. If you create the keys of node $n$ with a right neighbor $n$ so its left neighbor ($n$ so step 2 or step 3). If another keys to left ($n$ so step 2 or step 3), others a processor $n$ ($n$) has available capacity, when a processor $Q$ has available capacity, then transfer its own left, sets the right. If a previous can take all the keys, when a process — others, when a processor can take all the excess keys $K$, others, when a processor can take all the excess keys. This is equivalent in adding extent fails again it can just enter movement permanent.

- **Step 5.** See if $R$ is either the right neighbor's owner or the left neighbor's owner. If so, swap with its right neighbor or left neighbor by transferring the excess keys and stop.

- **Step 6.** Split the node $s$. Give the new sibling to processor $R$. Stop.

**Apparatus Merge.** In the above merge algorithm, we search through the set of nodes owned by the processor for one such that merge-candidate can take off keys efficiently in the aggregate range spread. As this search for a node that the processor owns still has a full range of keys, we need to compute a neighbor. If all of its right or left neighbors are full, then we compute one of them. We need a merge of keys to process keys. If there is a neighbor that is not full, then we use a merge of keys to process keys. If there is ...

- **Step 1.** For any leaf node $s$ to be (as a record by the processor $P$) that has excess keys $s$. Let the right neighbor of $s$ be $r$ and the left neighbor be $l$. Merge with the neighbor by transferring the maximum number of keys if and only if it can.

- **Step 2.** Scan through the list and add this new code to the list as long as it represents the maximum of $s$ and $r$.

- **Step 3.** Pick the last node on the list as a real and by the processor $P$ then has less than the minimum number of keys, let the right neighbor be $r$. Now if there is a neighbor that is full, then we use a merge of keys to process keys. If there is ...

- **Step 4.** For any leaf node $s$ to be (as a record by the processor $P$) that has excess keys $s$. Let the right neighbor of $s$ be $r$ and the left neighbor be $l$. Merge with the neighbor by transferring the maximum number of keys if and only if it can.

**Step 5.** If a minimum of $k$, then go to step 6 of the merge algorithm. The merge packs give the node that can be merged with its right neighbor by giving away its minimum keys. Stop.

## Experiments, Results and Discussion

The cloudiness of the different studies is similar to that of the AR test, except that the lower field has targets (systems) and can hold an arbitrary number of keys. The system studies have an average funnel, defined as 56% of the consistent funnel. We performed experiments with average funnels of 4, 7 and 10. A total of 700,000 keys were associated and upto 20 persons were used for distributing the B run.

## Experimental Setup.

A uniform random distribution of keys is chosen to create the actual AR test. Initially each presence was given one leaf node with a range of keys.

To study the load variation behavior under overation, we collected distributed snapshot of the processes at intervals of every 10,000 keys inserted in the lifetime. At each snapshot, we record the processor capacity in terms of the number of keys, the number of nodes in the load, the number of the number of node splits, merges and deletions. We also record the number of queries a presence receives, the number of messages it sends, and the amount of data balancing algorithm and the number of nodes that a transfer.

Section in the AR test, other experiment variables are the number of message keys for a search, the width of explanation and the number of nodes required by a transfer caused by a load balancing. We also calculated the amount of nodes that a node transfers before attaining a stable state when the number of nodes in the system increases.

**Results.** We first compared the version and the merge algorithms. In this experiment we built a AR test of an average funnel of 10 at the uniform nodes, with 800,000 keys and over load. At the 20 processes. We observed the two algorithms behaved quite similarly for certain statistics. We noticed both the algorithms did a good job at maintaining a data balance with the same being around 7% and the

variance being 0.00001). The number of hops per attempt also rises similarly in both algorithms, from 1.16 to 2.46 which surpasses the maximum hops (this filters a tree of height 3. The width of confidence varied between 4.6 to 7.0 for the algorithm



Figure 5.12: 40-site: Comparison of the Random vs Merge Algorithms

The difference in the algorithms is reflected in the number of hosts and the number of retries under that are stored at each generator. We see from the graph 5.40c that a 40-tree decreased over 50 generators the random algorithm stores 5000 hosts whereas the merge algorithm stores around 1000 hosts. This shows that the merge algorithm uses a far superior job at reducing the storage overhead of the 40-tree. However, the number of merges that occur is about 1000 for the random algorithm whereas for the merge algorithm this number is 2000. Also, since as a retries stay constant for the merge algorithm with 20 nodes and 200 expiries being tracked (c.f. section) on the retransmissing while only 26 nodes and 72 expiries are tracked with the random algorithm.

The results obtained from the above algorithms indicate definitely that the merge algorithm is more efficient at reducing storage space without affecting the number of hops per attempt and the width of confidence. So we decided to explore the

Table 3.6 Merge Algorithm. Comparison of different sizes with 3.5 Million and 5 Million Keys

| Keys | Processors | Leaves | Internal Nodes | Vector Searched | Object Searched | Splits | Merges |
|------|-----------|--------|----------------|-----------------|-----------------|--------|--------|
| 3.5 Million | 10 | 160 | 14 | 113 | 1320 | 517 | 2271 |
| | 20 | 155 | 10 | 136 | 964 | 566 | 2411 |
| | 50 | 161 | 8 | 152 | 616 | 695 | 2601 |
| | 80 | 151 | 11 | 154 | 528 | 801 | 2500 |
| | 100 | 163 | 10 | 172 | 531 | 2300 | 4955 |
| 5 Million | 10 | 310 | 15 | 168 | 2420 | 511 | 3471 |
| | 20 | 302 | 10 | 180 | 1582 | 565 | 4361 |
| | 50 | 300 | 8 | 192 | 980 | 607 | 3615 |
| | 80 | 301 | 10 | 196 | 682 | 610 | 3635 |
| | 100 | 303 | 14 | 186 | 534 | 1533 | 4359 |

One of the aims of this work was to see the effect of the number of keys or the effect on the number of processors and internal nodes. So we performed experiments with 3.5 million and 5 million keys.

From the Table 3.6 we see that for 10 processors, on a different with 3.5 million keys, the number of leaves is 160, internal nodes is 14, the number of vectors searched for restructuring is 113 and object 1320. The corresponding numbers for a different with 5 million keys is 310, internal nodes 15, vectors 168, and object 2420. The number of splits and merges vary depending on the number of processors. We also see the number of processors increase the number of merges also increases. From Figure 3.12 and 3.13, it is also seen that the number of splits and merges are not greatly proportional to...

It can also be seen that both the number of leaves and internal nodes increase nearly linearly with the number of processors (Figures 3.14, 3.15, 3.16 and 3.17)

Figure 5.44  Effect of increasing the Number of Processors on the Number of Leaves stored in a iR2-tree with 3.5 million Keys for the Merge Algorithm



Figure 5.45  Effect of increasing the Number of Processors on the Number of Inverse Nodes stored in a iR2-tree with 3.5 million Keys for the Merge Algorithm

Figure 5.16  Effect of Increasing the Number of Processors on the Number of Leaves stored in a 30 Leaf with 5 million Keys for the Merge Algorithm.



Figure 5.17  Effect of Increasing the Number of Processors on the Number of Ignition Nodes stored in a 15 Leaf with 5 million Keys for the Merge Algorithm.

Another observation is that the number of repeat touched by males/caring/grooms is less constant as the number of caretakers increases.

Thus far, we observed that the newest algorithm reduces the number of leaves and external outliers, and increasing the size of the 40 rate does not have a great effect on the number of leaves and victims for those. On the next step we to see if we could reduce the number of leaves even further by varying more input parameters in the experiment.

**Extensions:** The main objective of performing further experiments is to observe the effect of the main parameters on the number of leaves at the 40 rate.

- Questions to be answered:

  - What are the input parameters that affect the outputs of the 40 rate?

  - How does one vary the selected input parameters?

- Experiment. In our original experiment, we created 0.5 million legs in total at a 40 rate with an average funnel of 10. We built one initial 40 rate by assigning some number of legs to each caretaker. After the initial 40 rate is built, legs can control the latest rate to grow. When a caretaker decides that it holds too many legs it invokes the leaf funnel, which attempts to distribute the legs among the entire processes. If none of the entire processes have available capacity, then a caretaker is frozen and its capacity increased by an increment. We cannot find by selecting the number of legs to build the initial 40 rate and the funnel rate appropriately, we could reduce the number of victims in the final 40 rate. Thus, we chose the following three to build a small initial 40 rate, and the increment which is the average added in a processes during leaf

Table 5.1: Comparison of Doubling Initial Keys and Increments for a 4K-size with 3.3 Million Keys

| Algo | Processors | Locality | Interior | Border | Corner | Split | Merge |
|---|---|---|---|---|---|---|---|
| Doubling | 20 | 128 | 159 | 140 | 53 | 1058 | 188 |
| Initial Keys | 10 | 324 | 148 | 152 | 706 | 1643 | 507 |
| Original | 20 | 460 | 353 | 84 | 935 | 3201 | 3093 |
| Increment | 40 | 1073 | 947 | 141 | 1533 | 5305 | 5141 |
| Ground | 20 | 216 | 35 | 97 | 460 | 555 | 166 |
| Initial Keys | 10 | 300 | 78 | 97 | 760 | 856 | 846 |
| Border | 20 | 1724 | 139 | 141 | 842 | 2909 | 2744 |
| Increment | 40 | 3344 | 341 | 76 | 1153 | 3018 | 3374 |

behaving, as the set organ parameters to vary. The increment changes the route for the entire run.

The next concern is how to vary these parameters. We started off by allowing a growth of 4k times for the 4K-size, hence if the final 4k-size holds 3.3 million keys, then doubled keys is shown as 5 hundred (256=number of processors). The increment is shown as 8 *control keys*.

■ Variant Increments and their Benefits

We have the following outcomes:

— *Double Initial Keys: Original Increment*
In this scenario we allowed a growth of the over 4k times, so we control $1/3$rd (first/256=number of processors)=1/3rd respectively. The increment was chosen as $1/473/256$=c=color/processor/. The number of keys increases

at the end of the run was 565 and the number of nascent males was 59. The number of males marked for maintenance was 116 and copies 104 [for 10 processors] (Table 3.7).

- *Original Initial.step, Double Increment*

Here, we allowed a speed-if 56 times for the run and hence nascent 6 broodloss(Micronuclei)processors and thatfold the necessary to 2+(2+2 broodloss(150 x numbers/processors)). The number of loses reared at the end of the run was 565 with nascent males being 59. Males marked for maintenance was 206 and copies 420 (Table 3.7).

The numbers obtained above show that it's initial,step runs, double increment method reduces the number of forces evenly by half. Comparing this to the original algorithm, with initial,step = 1.0(150 x numbers/processors) and increment = 2 x 2/(150 x numbers/processors) we see that the number of forces has reduced from 104 to 52 and cortone codes from 104 to 83.

The results obtained from performing the experiments were something enough to permit us to reduce the effect of the reduction of fold,step increment. Notably! We noticed that it was the increment that was added to a given time was called to the number of forces in the set called to the *Micro*-. Hence, we varied the increment keeping the Original Initial.step.

After performing the experiments with these two increment to investigate for time we came up with three other measures. We added the following variables (as listed in table 3.8):

- *Original Initial.step, Half Increment*

- *Original Initial.step, Quarter Increment*

Table 3.8  Variance Scenarios of the Input Parameters for a 10 Gram of 2.5 Million Keys.

| Scenario | Number of Keys | Covariant |
|---|---|---|
| Original Keys Original Increment | $K$ | 1 |
| Original Keys Double Increment | $K$ | 2*1 |
| Keys doubled Original Increment | $2*K$ | 1 |
| Original Keys Half Increment | $K$ | 1/2 |
| Original Keys Quarter Increment | $K$ | 1/4 |
| Original Keys Triple (3 counts of processes) | $K$ | 3/10 |

$K = 2.5$ M (number of processes)

$1 = 2 * K$

Table 5.9: Effect of Changing the Increment on a dE List with 2.5 Million Keys

| Increment | Processors | Leaves | Internal Nodes | Leaves Leaves | Internal Leaves | Copies | Splits |
|---|---|---|---|---|---|---|---|
| IP Increment | 56 | 283 | 31 | 41 | 440 | 388 | 1038 |
| Increment1 | 56 | 398 | 141 | 125 | 159 | 534 | 567 |
| Increment1 | 16 | 523 | 161 | 196 | 208 | 1012 | 8575 |
| Increment4 | 8 | 645 | 302 | 298 | 301 | 2151 | 19102 |
| Increment10 | 16 | 6668 | 542 | 178 | 1364 | 16007 | 60732 |

*— Disposal Initial Step: Yield (Iterations)*

We observed that halving the increment increased the number of leaves to 50%, reducing the increment to a quarter of its original brought the number of leaves to 56% and storage nodes to 20%. Then, changing the increment from 3 to 1 and then to 2, changed the number of leaves from 283 to 293 to 293 (table 5.9), showing an almost linear dependence of the number of leaves on the increment added to a positional during load balancing. Thus, like we can conclude that the size of the tree depends on the number of times the increment is performed. With a small increment the number of times the increment is performed might lead to the increase of the tree and also, the creation of leaves.

The above discussion has shown that the merge algorithm and its success are more closely related to how the increment used for the number of times the number of leaves might happen, we could improve leaves by designing a new algorithm, or we developed the aggressive merge algorithm. We chose a new variant of the aggressive merge algorithm to operate on the aggressive merge algorithm, or figure. I also has

48 processors. The number of leaves for the merge algorithm is about 3940, whereas for the aggressive merge the number of leaves is only 320. The plot shown is that the aggressive merge algorithm is definitely more efficient.



Figure 5.48 :Course Comparison of the Merge vs Aggressive Merge Algorithm

We also plot the number of leaves versus processor in Figures 5.49a and 5.49b and note the quadratic behavior of the curves. So, the aggressive merge algorithm does a much better job at reducing the storage overhead at each processor, while increasing the cost of restructuring as expected.

Figure 5.18  4E-line: Number of Leaves versus Keys for 58 puicessors for Agglomerate Merge Algorithm



Figure 5.19  4E-line: Number of Leaves versus Keys for 58 processors for Agglomerate Merge Algorithm

Figure 5.31 4D-test: Number of Leaves versus Keys for 24 processors for Aggressive Merge Algorithm



Figure 5.32 4D-test: Number of Leaves versus Keys for 48 processors for Aggressive Merge Algorithm

Figure 5.22. *k*5.save: Number of Leaves versus Keys for 16 processors for Appleseen Merge Algorithm



Figure 5.24. *k*6.save: Number of Leaves versus Processors for Appleseen Merge Algorithm

In figures 5.49 through 5.51, we plot the number of leaves versus the number of legs for different numbers of processors, varying three between 12 and 20, for the aggressive setup algorithm. We also plot the number of leaves versus processors for 5 million legs and note the quadratic curves of the curves (Figure 5.50). It can be seen from the charts 5.49 and 5.50) that the number of leaves is balancing well, reaching a plateau for the plot of 50 and 20 processors. A good algorithm should have no more than about six 12/6 leaf nodes (a processor to complete with every other one). Our aggressive merge algorithm achieves this as the number of leaves flatten out with increasing numbers of legs for 50 and 20 processors. For 20 million and more processors the simulation did not execute long enough to reach a plateau value, as the final number of leaves is less than six 12/6 for n = 50.

As for the off-line algorithm, here we are observed the mode of replicators at all levels, the height of the tree and the number of legs per message for a off-tree with 5 million legs. We see that the height of the off-tree is 1 for 15 processors and 6 for 50 to 20 processors, with the number of legs varying from 1/25 for 15 as we increase the processors from 12 to 20. The width of replicators at level 2 never between 6 and 2 to 8/10. We then see that our algorithm does not significantly increase the tree and message overhead at this level.

### Summary

This chapter has thus far concentrated on the performance of replicators and balancing algorithms from a qualitative point of view, by doing large scale simulations. There are we measured only characteristics such as tree latency, message route and throughput. The timing information gives us an idea of how fast the system can respond and throughput. The timing optimization gives us an idea of how fast the system

responds to a query and what the throughput of the system is. A view of the state law of queues at one person per second. To obtain these timings, we go back to the implementation of our distributed E-tree that we discussed in Chapter 6.

### 7.4.1 System Response Time

Response time means the time taken for a single query to be processed. The author process work on a operation query and work fill an answer issues back from a node messages that the operation has been completed. The total time taken for the operation to complete is noted. The author then work out each per query. The average of the time taken for all operations gives the response time for a single operation. Response time is defined as

response time = total time taken for all queries / number of queries

Operation rate is defined as

procedure rate = total number of concepts processed / time taken to process the resources



Figure 6.10: Experimental Model for Measuring System Throughput

155

**Experiment.** Each generator is a *capitalist private strategy* rules that generates the messages ([2 part 3.3]). The possession of the messages is generated by the generator. In each generator, a sender, a wallet, and their possession which spawns the messages for a given set of time that make the message for a certain instant for the subsidy, the possession. Each message transmits to the generator at the time the wallet for the generator in a simplified, the node occupies time makes the message and returns it to the author. The author then obtains the time index for each message. After it receives all the messages, it then calculates the average response time and generates a set.

We have chosen to discern the response rate of $A$, $B$, and it processes. In our experiments, we obtained different generators rates by varying the clock interval between successive messages and asked for response times. Our experiments show that the response time decreases as the generators rate increases. This could be because of lesser time order in the file. Some may occur earlier earlier than others for a variant generators rate the response time is expected to increase as the system of clocks we find that there are fewer duties too in the system, and so the node treats the file as if it were unimportant, over though we replaced with generators that happened at the same time. This response time is inversely proportional to the propagation rate for all generators with the same propagation time. This can open up as many possible (with as deep interval) a not sufficiently large enough to block the system will pause. So all the generators in every case is different sufficiently large to the end at which case we do not consider the generators time as the system response.

**Results.** The graphs show the response times and generation rates for $A$, $B$ and $C$ processes.

Figure 5.9: Response Times for a 4 Processor System



Figure 5.10: Response Times for a 8 Processor System

It was observed that for a 4 processor system, it takes about 35 milliseconds for an operation to complete, for a 8 processor version it takes 40 milliseconds and for an 8 processor system, the response time is 55 milliseconds.

In order to justify these findings, it is interesting to know the message transit times, processing times of a processor and queuing time. It is difficult to get an estimate of the queuing time, but we performed some sample experiments to determine the message transit time and the processing time.

Figure 5.10: Response Times for a 4 Processor System

**Degradation's Model.** For the processing time of any processor—task processing times, we use the same model that we used in earlier attempts. When the node manager receives a message, it now returns it with the processing start time and after processing the message, it again it can occupy the message with the processing end time. All the messages are returned to the anchor, so the anchor calculates the average total processing time.

To calculate the average total it has between any two processors and message time, the anchor approach of processor on different machines and messages are sent back and forth between all the processors and the anchor. Each message is time-stamped with the time that it was sent and the time it was received at another processor. The anchor then finally collects all the messages and calculates the average time it takes for a message to travel between any two processors.

Since our experiments we observed that the average travel time between two processors is approximately 3.2 milliseconds and the total processing time is around 4.8 milliseconds in table 5.12 we calculate the processing times and average travel times.

Table 4.18 Timing Calculations

| Parameter | Number of Slices | Preceding Time | Message Time | Total Time | User Calculations |
|---|---|---|---|---|---|
| 1 | 1.5 | $(1+t) * t$<br>$= 11.98$ | $8 * t$<br>$+10.72$<br>$= 8.78$ | 20.76 | $30.00.15$<br>$0.51$ |
| 1 | 2.5 | $(1+t) * t$<br>$= 12.73$ | $(10 * t$<br>$+10.72$<br>$= 12.98$ | 25.94 | 20.00.03<br>$0.35$ |
| 2 | 2.5 | $(2+t) * t$<br>$= 12.73$ | $(10 * t$<br>$+10.72$<br>$= 12.98$ | 24.96 | 20.00.05<br>$0.37$ |

Unit processing time is 5.1 milliseconds.
Unit transfer time is 1.5 milliseconds.

times as follows:

    processing time = (number of hops + 1) * unit processing time.

and the message time is

    message time = (number of hops) * unit message time + overhead time to initiate.

The checking time in nodes is added when a message starts at the author and returns to the author.

The time difference in the table 4.18 can be attributed to the delays that outside message collisions produce between the delays without delay into the ...

## 4.3. Performance Model

In this section, we present a simple analytical model that permits one to answer this and the maximum throughput of the distributed search structure we studied in this paper. The performance depends on the structure of the file tree or DB tree. For example, both the number of tree per increment and the degree of ...

replication affect the amount of overhead required to maintain the match structure. These values are very difficult to calculate, and they depend on the algorithm used to perform the data balancing. For this reason, we will use the estimate of the number of hops and the degree of reduction developed in Sections 3.3.1. The model described in this section is loosely based on the model presented in [10]. We assume that operations are generated uniformly to all processors, and the amount are made to the data uniformly.

We first define the variables that we use in the analysis:

- $S$ Number of levels in the match structure (level 0 is the leaf, level $L$ is the root).
- $P$ Number of processors that maintain the match structure.
- $R$ Average number of hops required to navigate to a leaf.
- $B_i$ Degree of replication at level $i$, $i = 1, \ldots, L$. $B_0 = 1$ and $B_L = P$.
- $C$ Maximum cycle boost.
- $n_i$ Probability that an operation is an insert operation.
- $p_{nq}$ Probability that an operation causes redistributing (join or merge).
- $t_n$ Time to process an action.
- $t_m$ Time to process a message.
- $t_a$ Processing time for sending and receiving a message.
- $t$ Arrival rate of operations to an individual search structure.
- $t_{op}$ Total arrival rate of operations to the distributed search structure.
- $N_o$ Average number of actions generated by an operation.

$N_w$   Average number of messages provided by an operation

$W$   Waiting time

$T$   Response time of an operation

$Th_{max}$   Maximum throughput

We start by determining the number of messages and actions required to process an operation, $R_s$ and $N_w$. Since there are $k$ levels, $L$ search actions are required from superior states (hops; if $L$ messages are required to amplify an operation/action). In addition, an operation needs state transitions for update, and all copies of the parent must be relocated about the new sibling. At once the parent-right split, with probability $p_{ins}$. Therefore,

$$R_s = L + k \sum_{i=1}^{L-1} (2M_i + R_{w,i})$$ (5.1)

$$N_w = R + \frac{1}{k} \sum_{i=1}^{L-1} (2M_i + R_{w,i} - 1) \times r$$ (5.2)

If $k$ is the rate at which operations are generated at a node that helps us maintain the distributed search structure, then the total rate at which operations are generated is

$$\lambda_{op} = Pk$$ (5.3)

A processor that helps to maintain the distributed search structure will be required to process this compound operation in addition to ordinary search and join jobs that are not connected to average passing. The average time to process a job is

$$\lambda_{tot} = (N_w A_s + N_w A_w)(N_w + N_{op})$$ (5.4)

Since the cost is fully replicated, it is not a bottleneck. If the data balancing distributes the nodes properly, then no leaf node is a bottleneck either. Therefore, the work to execute an operation is evenly spread among the processors in the system.

As a result, the processor utilization due to search structure processing is:

$$p = \lambda \left( N_s t_s + N_L t_L \right) \tag{5.1}$$

The time that a job spends waiting for processor service can now be calculated by applying a queuing model. We use a simple $M/M/1$ queue, and find that:

$$W = t_w \frac{p}{1-p} \tag{5.2}$$

The time to get a response from an operation is the time to remove all messages and actions associated with the operation:

$$T = 2(W + t_N) + t_s + t_L + (Q(W + t_s + t_L)) \tag{5.3}$$

The maximum throughput is the maximum rate at which every processor can execute the jobs associated with the search structure operations:

$$Th_{max} = N/(N_s t_s + N_L t_L) \tag{5.4}$$

In a distributed search structure with a large number of processors, the overhead of monitoring the search structure is generally that in the number of legs, $H$, and the cost of monitoring the level $S$ nodes. As we saw in Section 5.3.1, $H$ approaches an asymptote for a fixed height case. The algorithms described in [20] ensure $H_s$ arrives for every split of a level 1 node. Fortunately, we found that $H_s$ grows very slowly with increasing $N$. As a result, the overhead of monitoring a different does not increase in size in the system increases when processors are added. As such, the $SD$ tree algorithm is scalable in a very large number of processors.

## 5.5.1 An Auditorium

- Analyze for a 8 processor ES tree.

Let us make an analysis of a silicate-rhodochrosite tree in paramount for which we have performed the kriging experiments. So, we have $N = 8$ and average lumen $f = 50$. In Section 5.5.1, we saw that in a large format silicate role 1 batch, the number of keys is about 1, and the well-of-emphasis up level 1 is about $1\,500 = 4000 \times 8$, where $P$ is the number of paramount. We have found that the level 2 nodes are equivalent at nearly half the number of paramounts, so we will assume that $R_0 = P25$. We measured the time to process a message at $t_c = 30$ns seconds and transmission time for a message in $t_n = 3000$ seconds.

(5.10)

With these statistics in mind, we will use the following additional parameters as input to this model:

$$t_m = \qquad 881$$
$$k_c = \qquad 1$$
$$p_{max} = \quad 1/(f) = 1$$

We use these parameters to determine the number of messages and streams that an injection generates:

$$N_i = 1.8k3$$
$$N_m = 3.0k8$$

We can use the determination of the number of actions and passages to compute the average correction rates and the streams interface throughput:

$$t_{c,p} = 4k2B$$

$$\overline{x}_{n,m} = 343$$

With a processing rate of 331 operations per second $\mu = 1/\overline{x}$ and the response time for an operation is 0.003 seconds. From the chart 3.16, we see that the lowest response time is around 900 seconds. The analysis gives a more pessimistic value taking into account some running time.

• Analysis for a 3d processor 4-bit bus

We will do a similar analysis for larger dB sizes that we used in our experiments in section 3.6.1 with $J = 56$ and average clusters $J = 33$. We see that, as in the case of figure 3.1, at level 0, $d_0 = 10$, at level 1, $N_1 = 256$ and at level 2, $d_2 = 10$ at level 3, $N_3 = 1024$. The number of keys is $J = 56$.

Again, similar to the analysis for 4 processors, here too we use additional parameters as input to the model:

| $I_n =$ | 8044 |
| $k_n =$ | 3025 |
| $L_n =$ | 886 |
| $q =$ | 5 |
| $d_{n,1} =$ | $1/(J = 343$ |

We then compute the number of messages and arrive thus at question generation

$$N_q = 1458$$

$$N_m = 11315$$

We now calculate the average expression time and maximum throughput.

$$t_{avg} = .0059$$

$$T_{R_{max}} = .0417$$

With a processing rate of 1268 operations per second, $\rho = .172$, and the response time for an operation at 1268 seconds.

For a comparison, consider the performance of a centralized index server that has the same message passing costs, $t_m = .001$. Servicing each request requires the processing of ten messages (the request and the response). We will assume that the actual index lookup requires $t_v = .0044$ seconds. Then, servicing an operation requires .0064 seconds, allowing a maximum throughput of 156 operations per second, 15 percent reduction in throughput. Conversely, the response time for an operation is .026 seconds. Therefore, at the cost of doubled latency, the throughput is increased by a factor of 18 by using the distributed search structure.

## 5.6   Conclusions

In this chapter, we have described extensively all the algorithms developed, experiments conducted and the performance results obtained. Here, a brief summary is presented by listing the conclusions drawn from our experiments.

• Replication: In section 5.2 we have presented two algorithms for replication namely, *full replication* and *path replication* and discussed the method of maintaining indices coherency. To compare the performance of the algorithms we examined the overhead of the full and path replication of the active nodes and found that path replication imposes much less overhead in terms of space and

messages that fall replication. The width of replication measure be path-replication shows a reliance measure with number of processors and hence provides a suitable short-term index.

- **Data Balancing.** We also conducted simulations on a large scale to calculate the results obtained from our implementation. The simulation results show that our algorithms for data load-balancing schemes give much better results than other proposed schemes. In our simulated work, we attempt to show that load-balancing overhead is not high. The simulated and distributed data balances indicate performance overhead is well suited for realistic situations.

We observed that:

- **Incremental Growth/Data Balancing.** The results of this experiment are similar to those of the general algorithms with the width of replication being 1.5 and the number of hops around 3.6 for 38 processors.

- **Fixed-Replication-Balancing.** We performed experiments with fixed-degree trees of 3, 5, and 7, with features varying from 10 to 50. The width of all these increases almost a platau with increasing bound. The value of the width of replication ranges and between 3.6 and the number of hops is 3.45. This is in accordance with the formula we have derived for the width of replication at level 2 which is

width of replication at level 2 = 1.960 × .802$^{3}$, where T is the number of processes.

We also notice that the width of replication and number of bays depend only on the number of processes. The fixed length AB now represents a flow-line that repetitions does not depend on the larger ones with a larger bound. Thus all the algorithms analysis through scheme.

- *B case.* We also designed the individual aware case that is then balanced on the sub number of bays with large σ... process. We first computed two algorithms, varians and array. Of those we found that the array algorithm determines the number of ordinary aware entries in a parameter when it runs the individual as by process. Larger the individual away and its numerons orbital in a parameter when it runs aware. We hoped that the number of ordinary varies linearly with the position of vertices; this means ever the number of main aware in summented is achieved, and the number of bays is the same as in the AB case. We then developed the aggregate merge algorithm, where we while the mowing so many bays as needed. We hasn't line of the view the aggregate merge in the A-bit, whereby the space overhead latest of a bays... process, indicating the number by multi ratio = 1/2, which is typically much simpler less the straight of T process. As a answer, we expect its replication to a measure. The enexpectation (x) is a summary.

  - A difficulty that is occasioned when the most... of the individual aware is between max x − 1/2, which is equivalent to replication to the... A from this view the aware factory a then is between max x − 1/2, which is equivalent to cuts overhead bays. We then developed to a measure, are thing only a small more x. Using table 5.12 to account for the findings we have calculated indirectly. Using table 5.12 to account for the findings we have calculated

- *Testing Study.* We performed testing experiments on our implementation in a pather same idea as the system response timer for result and insert parts were. We find that, with 1 processors the system resource timer was one... all the individuals implements. Using table 5.12 to account for the findings we have calculated

- Analytical Performance Model : We used the characteristics of the large scale dB rise to develop a sample analytical performance model. We then studied the effect of increasing the number of processors, and found the potential of maximizing the dB rate grows very slowly. We applied the performance model to analyze the results obtained from our experiments on the dB rise with 8 processors and 16 processors. With both analyses we found that the model provides a slightly larger response time than what we obtained by our experiments. Our experiments provided a response time of 56 milliseconds(figure: 5.18) while the model predicted 58 milliseconds. Finally, from the model the observation was that the overhead of maintaining a dB tree is not significantly affected by the inserts factor, as long as the factor is large. We found that a distributed search structure performs much longer throughout than a centralized index access, at the cost of a modestly increased response time.

## CHAPTER 6
## CONCLUSIONS

In this dissertation we have worked on distributed B-trees. Our contributions in this fact form the development and implementation of several algorithms for their balancing a distributed algorithm B-trees. In Chapter 1 we presented the goal of our work and provided the motivation for pursuing this research. We also presented some background on distributed data structures. We outlined the B-tree version of our feasibility as a distributed structure.

In Chapter 2 we discussed concurrent B-trees and distributed B-trees. We also presented useful applications of the distributed B-tree, namely the distributed indexed tree, the off-line and on candidates for parallel merged file systems.

In Chapter 3 we presented the theoretical framework of the replication algorithms developed by us. We presented two approaches, lazy synthesis copies and variable copies. We also synthesized these algorithms and they are termed full replication and path replication, by the purposes of our implementations.

Chapter 4 presents the details of our implementations; the underlying architecture and details on the node migration mechanism that is fundamental to data balancing. We also presented the experiments protocol that is relevant to our data balancing algorithms. Other details, the node structure, node routing and updates are also discussed. We have also studied the portability of our implementations by porting it to the AIM, a shared memory multiprocessor system with 16 processors.

156

Finally, in Chapter 5 we presented all the algorithms that we have developed, for replication and data clustering. We discussed the algorithms and their performance in detail.

The performance results of the replication algorithms show that among the two methods of replication, pull-replication performs better and is suitable for scaling in large cases. It moves much less content than in full replication and hence is better for large datasets. The truth of replication does not increase linearly with the number of customers and hence it scales for huge datasets as well. When the latency is very high, we experience a worse performance but the benefit of replication is still apparent.

We developed centralized and distributed algorithms for data balancing and observed that distributed algorithms with sequential probing perform very well compared to the others. In terms of the total system load imbalance, the distributed algorithms with sequential probing perform as good as the centralized algorithms. The distributed algorithms take a much smaller amount of time to run since they execute in a decentralized manner. The parallel execution of the distributed algorithms makes them faster. In terms of the different amounts of incremental growth data, balancing and their result time also increases. The incremental growth performance shows consistent behaviour across different probing techniques, which quickly yields a better migration plan. We observed that the distributed algorithms with sequential probing perform very well and yield a result around 1.2 and the number of hops around 3.6 for 30 processors. We performed experiments for fixed height items of 3, 5 and 8 and found using means 30, 40 and 50 processors for growth ratio of 30% and 50%. We observed that the width and imbalance varies with varying items as for 30. We discussed the behaviour of these algorithms in detail.

We concluded the duration of our text, the different and performant data balancing algorithms. We discussed their algorithms for balancing, search, random, and aggressive range. Of these the aggressive range algorithm does the best in terms of reducing

a small data balance with negligible overhead. The asymptotic number of leaves in a *db* run using the aggressive merge algorithm is about $p(p-1)/2$, which is typically much smaller than the number of leaves in a *sfd* run. With the merge algorithm, we studied various extensions by changing some input parameters to the *db* tree and noted that the *db* tree performance was greatly affected by the increment size that is used to add storage to a processor when a run starts.

In order to determine how well our implementations works, we performed timing experiments and studied the response times of our system. With 8 processors we obtained a response time of 54 milliseconds. We have also provided an explanation of the storage we obtained (table A-4).

Lastly, we presented an analytical model to validate our experimental results. We applied the analytical model to analyze our experimental results and found that the model predicts a more profoundic time than the timing we obtained.

# REFERENCES

[1] Ahuja, S. Carlson T. and Gelernter, D., *Linda and Friends*, Computer, August 1986, Vol. 19. No. 8, pp. 26-34.

[2] Bal, H. E., and Tanenbaum, A. S., *Distributed Programming with Shared Data*, Proceedings of IEEE International Studies on Computer Languages, March 1988, St. Louis, MO.

[3] Bal, H. E., and Tanenbaum, A. S., *Distributed Programming with Shared Data*, Computer Languages, 1991, Vol. 16. No. 2, pp. 129-146.

[4] Bal, H. E., Kaashoek, M. F. and Tanenbaum, A. S., *Orca: A Language for Parallel Programming of Distributed Systems*, IEEE Transactions on Software Engineering, 1992, Vol. 18, No. 3, pp. 190-205.

[5] Baxter Idem, L., *Expected Behavior of B³-trees Under Random Inserts*, Acta Informatica, 1982, Vol. 19. No. 1, pp. 158-173.

[6] Bensten, P. B. Spinger, V. J. and J Liou Vtr., *Concurrent Maintenance of Data Structures in a Real-World Environment*, The Computer Journal, 1990, Vol. 33, No. 3, pp. 155-170.

[7] Bayer, R. and McCreight, E., *Concurrency of Operations on B-trees*, Acta Informatica, 1972, Vol. 1, pp. 173-189.

[8] Bernstein, P. A. Hadzilacos, V. and Goodman, N., *Concurrency Control and Recovery in Database Systems*, Addison-Wesley Publishing Company, 1987.

[9] Birrell, A., *Operation Specific Locking in the Server*, Proceedings of the Principles of Database Systems, ACM SIGACT/SIGMOD, 1987.

[10] Ceronru, A. Golomon D., Matwin, C. J. and Shronko, B. P. *Linda Alternatives for Managing parallel Systems*, Parallel Hash, 1988, Vol. 18, No. 5, pp. 505-601.

[11] Clony, J. C. and Clay, L. N., *A Note on Measuring to Reliability Key-Packing of the Luxury Network Table*, (abstract on Systems 90), Vol. 37, No. 1, 1, pp. 60-64.

[12] Colbrook, A. Spirals E. K. Delecroiz C.S. and Wald, W. E. *An Algorithm for Concurrent Search Trees*, Proceedings of the 1991 International Conference in Parallel Processing, 1991, pp. 38-41.

[17] Dettelé,Berger, M. *How to Undertake a Dictionary as a Complete Network. Proceedings of the 22nd Annual ACM Symposium on Theory of Computing, Baltimore, MD, USA, May 1990, pp. 1-17. 189*

[18] Ellis, C.S. *Extensible Data Structures. A Case Study. IEEE International on Computers, 1985, Vol. c-34, No. 12, pp. 1178-1187.*

[19] Evangelist, M. *A Grammatical Justification Access Technique Mechine. IEEE Transactions on Parallel and Distributed Systems, 1991, Vol. 2, No. 1, pp. 148-152*

[20] Flachbar, V.   *Optimal Allocation of Branch Data Over Distributed Memory Mechanism. Proceedings of the 5th International Parallel Processing Symposium, Beverly Hills, CA, USA, March 1991, pp. 568-556*

[21] Dooley, M.   *A Methodology for Implementing Highly Concurrent Data Structures. Second GEM REPLAN Symposium on Principles and Practices of Parallel Programming, Seattle, WA, March 1990, pp. 197-206*

[22] Herlihy, M. and Wing, J. Linearizability. *A Correctness Condition for Concurrent Objects. ACM Transactions on Programming Languages and Systems, 1990, Vol. 12, No. 3, pp. 463-492*

[23] Herlihy, M. *Hybrid Concurrency Control for Abstract Data Types. Journal of Computer and System Sciences, 1990, Vol. 43, No. 1, pp. 25-61*

[24] Herlihy, M. *A Methodology for Implementing Highly Concurrent Data Objects. ACM Transactions on Programming Languages and Systems, 1993. Vol. 15, No. 5, pp. 745-770*

[25] Johnson, T. and Shasha, D. *Utilization of B-trees with Inserts, Deletes and Modifies. Proceedings of the 8th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, Philadelphia, PA, USA, March 1989, pp. 235-246*

[26] Johnson, T. and Shasha, D. *A Framework for the Performance Analysis of Concurrent B-tree Algorithms, Proceedings of the 9th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, Nashville, TN, USA, April 1990, pp. 273-287*

[27] Johnson, T. and Colbrook, A. *A Distributed Data-Balanced Dictionary Based on the Delah Tree, Proceedings of the 4th International Parallel Symposium Symposium, Beverly Hills, CA, USA, March 1992, pp. 319-324*

[28] Johnson, T. Krishna P. and Colbrook A. *Distributed Indices for Accessing Distributed Data, Proceedings of the 12th IEEE Symposium on Mass Storage Systems, Monterey, CA, USA, April 1993, pp. 199-207*

[29] Jones, D. *Concurrent Operations on Priority Queues. Communications of the ACM, 1989, Vol. 32, No. 1, pp. 132-137*

[30] Kung, H.T. and Lehman, P.L. *Concurrent Manipulation of Binary Search Trees. ACM Transactions on Database Systems, 1980, Vol. 5, No. 3, pp. 354-382.*

[31] Kumar, V. and Krishna, P. Deep Indices for Distributed Search Structures, Proceedings of the 10th ACM SIGMOD, International Conference on Management of Data, Washington, DC, USA, May 1993, pp. 321-341*

[17] Johnson, T., *Supporting Assertions and Deletions in Striped Parallel Filesystems*. Proceedings of the 5th International Parallel Processing Symposium, Newport, CA, USA, April 1992, pp. 125-132.

[18] Johnson, T. and Shasha, D., *The Performance of Concurrent Data Structure Algorithms: ACM Transactions on Database Systems, 1993. Vol. 18, No. 1, pp 51-101.

[19] Jul, E., Levy, H., Hutchinson, N. and Black, A., *Fine Grained Mobility in the Emerald System*, ACM Transactions on Computer Systems, 1988. Vol. 6, No. 1, pp. 109-133.

[20] Kjelleruff and Johnson T., *Supporting distributed search structures: Technical Report TR93-002. University of Florida, 1993.

[21] KMG: Principles of Operation, Kendall Research Corporation Copyright Publications, 1991.

[22] Kindberf, T. and McEntosch, P., *Page-ring Based Data Structures on Distributed Memory Architectures*, Journal ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming, St. Petersburg, FL, March 1993, pp. 139-148.

[23] Kuhn, D. and Ellis D. E., *Pivotal Publishing Techniques for Parallel Join Systems*, First International Conference on Parallel and Distributed Information Systems, Miami, FL, USA, 1991, pp. 102-110.

[24] Kumar, V. and Segev, A. *Cost and Availability Tradeoffs in Replicated Data Concurrency Control*, ACM Transactions on Database Systems. 1993 Vol. 18, No. 1, pp. 102-131.

[25] Kuchler, J. and Langendoen, A., *Merle-2 HDD: A Cooperative Approach to Parallelization*, Computer Systems Journal, 1992. Vol. 5, No. 1, pp. 89-99.

[26] Law, F., Chen, V. and Baldazzo, J. M., *ARIES: A Transaction Recovery Method Supporting Fine-Granularity Locking and Partial Rollbacks Using Write-Ahead Logging*. ACM Transactions on Database Systems, 1992. Vol. 17, No. 1, pp. 94-162.

[27] Kuchler, J. and Langendoen, A., *Method for the Locking of Concurrent Operations in the Artikel-Server Parallelization Techniques and Parallel Processing Proceedings of Distributed Computing Systems*, Quebec, Canada, 1988, pp. 83-57.

[28] Kuchler, J. and Langendoen, A., *Merle-2 HDD: A Cooperative Approach to Parallelization*, Progress of Systems Journal, 1992. Vol. 5, No. 1, pp. 89-90.

[29] Law, F., Chen, V. and Baldazzo, J. M., *ARIES-1M: An Efficient and High Concurrency Index Management Method Using Write-Ahead Logging*, Proceedings of the ACM SIGMOD International Conference, San Diego, California, USA, June 1992, pp. 378-382.

[30] Leffand, S., *The User's Locking Protocol for B-Trees*, ACM Transactions on Database Systems, 1981. Vol. 6, No. 4, pp. 650-570.

[31] Lehmann, P. and Silberschatz A., *Distributed File Systems*, ACM Computing Surveys. Jourmal, ACM Computing Surveys, 1990, pp. 321-374.

[32] Lehman, P. L. and Yao, S. B., *Efficient Locking for Concurrent Operations on B-trees*, ACM Transactions on Database Systems, 1981. Vol. 6, No. 4, pp. 650-670.

[33] Lewin, M. A., Baminovska, R. Levy, G. and Kang, W. *Tree Pooling and Controlled Commit Block Migration for Multi-Tree Sequential Processing*, Vol. 5, No. 1, pp 108-913.

[16] Lukacs W., Setaro M., and Schneider, S. A., 24*. *Linear Scaling for the Solvation Free Proceedings of the 1993 ACM/IEEE International Conference on Management of Data*, Washington, D.c, USA, May 1993, pp. 207-216

[17] Lukacs W., Setaro M., and Schneider, S. A., 24*. *A Family of Order-Preserving Scalable Distributed Data Structures*, Proceedings of the fifth VLDB Conference, 1994, pp. 342-353

[18] Martikich, G. and Hommel, G., An Efficient Method for Distributing Search Structures, Proceedings of the First International Conference on Parallel and Distributed Information Systems (PDIS 1991), Miami, Florida, December 1991, pp. 108-117

[19] Miller K. and Snyder L., Multiple Access to B-trees Proceedings of the 1978 Conference on Information Sciences and Systems, Johns Hopkins University, Baltimore March 1978, pp. 400-406

[20] Fisher J. B., A Concurrent Search Structure, Journal of Parallel and Distributed Computing, 1988, Vol. 7, No. 2 pp. 186-291

[21] Peleg D., Distributed Data Structures: A Complexity Oriented View, Proceedings of the fourth International workshop on distributed Algorithms, Bari, Italy, September 1990, pp. 71-89

[22] Rao V. N. and Kumar V., Concurrent Access of Priority Queues, IEEE transactions on Computers, December 1988 Vol. 37, No. 12, pp. 1657-1665

[23] Rao V. N. and Kumar V., Superlinear Speedup From Parallel Search Concurrency, Proceedings of the International Conference on Parallel Processing, 1987

[24] Rosenberg, A. L. and Stockmeyer, L. and Snyder L, Uniform Data Encodings, Theoretical Computer Science, No. 1, March 1981, pp. 176-199

[25] Salzberg, B., Grid File Concurrency, Information Systems, 1986, Vol. 11, No. 3, pp. 235-244

[26] Salzberg, B. and Dimock A., Principles of Transaction-Based On-Line Reorganization, Proceedings of the 18th VLDB, Vancouver, Canada, August 1992, pp. 511-520

[27] Salzey, S. and Lomet D. Access Method Concurrency with Recovery, Proceedings of the 1992 ACM SIGMOD International Conference on Management of Data, San Diego, California, June 1992, pp. 351-360

[28] Sagiv Y., Concurrent Operations on B-trees with Overtaking, Journal of Computer and System Sciences, 1986 Vol. 33, No. 2, pp. 275-296

[29] Salzey, B. and Kaplan G., Disk Striping, International Conference on Data Engineering, Los Angeles, USA, February 1987, pp. 336-342

[30] Srinivasan, V. and Carey M., Performance of B-Tree Concurrency Control Algorithms, Proceedings of the 1991 ACM SIGMOD Conference, Denver, USA, May 29-31, 1991, pp. 416-426

[31] Sarger, V. and Lomet P. Multi-Disk Extent, Proceedings of the 2001 ACM SIGMOD, pages 138-151, 2001.

[32] Severance, C. and Pramanik, S. Distributed Linear Hashing for Main Memory Databases, International Conference on Parallel Processing, Illinois, August 1990, pp. 274-281

[33] Severance, C. Pramanik, S. and Wolberg P. Distributed Linear Hashing and Parallel Projection in Main Memory Databases. 16th International Conference on Very Large Data Bases: Brisbane Queensland, Australia, August 1990, pp. 674-682

[20] Shasha, D. and Goodman, N., Concurrent Search Structure Algorithms, ACM Transactions on Database Systems, 1988. Vol. 13, No. 1, pp. 53-90.

[21] Tang, J. and Natarajan, N., A Scheme for Increasing Availability in Partitioned Replicated Databases, Information Sciences, 1991, Vol. 53, No. 1-2, pp. 1-34.

[22] Ungureit, L., Bharghari V., and Weikum, G., Distributed File Organization with Scalable Cost/Performance, Proceedings of the 1990 ACM SIGMOD International Conference on Management of Data, Atlantic City, USA, May 23-25, 1990, pp. 253-264.

[23] Weihl, W. E., Commutativity-Based Concurrency Control for Abstract Data Types, IEEE Transactions on Computers, December 1988, Vol. 37, No. 12, pp. 1488-1505.

[24] Weihl, W. E. and Wang, P., Multi-Version Memory: Software Cache Management for Concurrent B-trees, Proceedings of the 2nd IEEE Symposium on Parallel and Distributed Processing, 1990, pp. 650-655.

[25] Yan, J. L. and Sadowe, E., Both Cells in Mutually Parallel Systems, Proceedings of the 6th International Parallel Processing Symposium, Beverly Hills, CA, USA, March 1992, pp. 380-389.

# BIOGRAPHICAL SKETCH

Padmasheree Lasor on Hyderabad, India is the second daughter of Mr. Achutu Aggarao and Mrs. Mameelalla. She has an elder sister, Lalithaji, a younger brother Ravi and a younger sister, Rajeshree.

She schooling was in Lasova Convent, Kurnli, a hill station in Bihar, India and her keen college education at St. Xaviers College, also at Kurnla. She then obtained a Master of Science in Physics from Central University of Hyderabad. She then proceeded to do her Master of Technology in Computer Science at the Indian Institute of Technology, Medras, India.

She was then executed as a Technical Officer in Electronics Corporation Of India Ltd., Hyderabad. She worked there for five years on Computer Graphics, Artificial Intelligence and Parallel Processing. She then decided to further her knowledge and snapped from the organization to pursue her Ph.D at the University of Florida.

Her current research interests include scientific databases and distributed systems.

Her hobbies include reading, jogging, cooking and handicrafts.

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Theodore J. Johnson, Chairman
Associate Professor of Computer and
Information Sciences

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Richard C. Newman-Wolfe
Associate Professor of Computer and
Information Sciences

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Randy Chow
Professor of Computer and Information Sciences

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Sartaj Sahni
Professor of Computer and Information Sciences

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.

Paul Avery
Associate Professor of Physics

This dissertation was submitted to the Graduate Faculty of the College of Engineering and to the Graduate School and was accepted as partial fulfillment of the requirements for the degree of Doctor of Philosophy.

May 1991

Winfred M. Phillips
Dean, College of Engineering

Louis A. Holbrook
Dean, Graduate School